



A potpourri of analysis considerations for microbial community function



Eric
Franzosa



Kelsey
Thompson

Curtis Huttenhower (chuttenh@hsph.harvard.edu)

2022-07-27



THE HARVARD CHAN
**MICROBIOME IN
PUBLIC HEALTH CENTER**



BROAD
INSTITUTE



HARVARD T.H. CHAN
SCHOOL OF PUBLIC HEALTH
Department of Biostatistics



Who am I and why am I talking to you?

- Background originally computer science and chemistry.
- Worked in the software industry for a few years.
- Got into computational biology during grad school.
- PhD in computational functional genomics / microbiology.
- Worked on microbial communities at Harvard since 2009 (!)
- Because they're cool!
And because they matter for population health.



The mission of the HCMPH

Use an understanding of the microbiome to improve population health, through basic research, translation, policy, education, and outreach.

The screenshot shows the website for The Harvard Chan Microbiome in Public Health Center. At the top, there's a red header bar with the Harvard T.H. Chan School of Public Health logo and a "Quicklinks" dropdown. Below the header, there are three circular icons representing different stages of life: a fetus, a baby, and a person. The main title "THE HARVARD CHAN MICROBIOME IN PUBLIC HEALTH CENTER" is displayed prominently. A navigation bar below the title includes links for Home, Research, People, Partnering, Resources, Contact, and HCMPH Symposium. The central content area features a circular phylogenetic tree diagram labeled "MICROBIOME". To the right of the tree, a quote by Michelle A. Williams is displayed: "*Understanding the microbiome may transform our understanding of how healthy bodies become diseased, how aging leads to infirmity, and how we could alter our internal ecosystems to prevent and treat a vast range of conditions.*" The quote is attributed to her as "Dean of the Faculty, Harvard T.H. Chan School of Public Health".



Curtis Huttenhower



Wendy Garrett



Andy Chan



Eric Rimm



Xochitl Morgan

<https://hcmph.sph.harvard.edu>

HCMPH activities

- **Research**
 - Large flagship consortium projects: OPTIMISTICC, HMBR, HMP, Micro-N...
 - Microbiome population and basic science studies
- **Facilities**
 - BIOM-Mass: BIObank for Microbiome research in Massachusetts (*and elsewhere*)
 - Harvard Chan Microbiome Collection and Analysis Cores
 - Harvard Chan Gnotobiotics Facility
- **Training**
 - IID209: The Human Microbiome and Microbial Communities
 - Year-round methodological short-courses and boot camps
 - Harvard Catalyst Introduction to 'Omics, Biomarker Science, Network Medicine
- **Community**
 - Microbiome Epidemiology Working Group (MEWG)
 - HCMPH monthly Scientific Meeting, annual Symposium
- **Collaboration**
 - Harvard Chan School, Longwood, Broad Institute, national, international...
 - Industry partnerships



THE HARVARD CHAN
MICROBIOME IN
PUBLIC HEALTH CENTER



Stuff we do in the lab

- About 85% computational:
 - ~Half methods development for microbial communities generally.
 - Computational
 - Statistics
 - ~Half applications.
 - Mostly human microbiome epidemiology.
 - Often gut, often inflammatory disease (IBD, arthritis, cancer, etc.)
- The rest is in the wet lab...
 - Data generation
 - Technology development
 - Validation experiments
- *How do microbial communities work and why?*



The bioBakery: a next-generation environment for microbiome analyses

<https://huttenhower.sph.harvard.edu/biobakery>

- Environment for meta'ome analysis
 - Shotgun metagenomes/metatranscriptomes, amplicons (16S/ITS)
 - Taxonomic and functional profiling
 - Experimental design, statistical analysis

- Pre-built one-click environment to run:
 - On your laptop graphically
 - On a server remotely
 - On the cloud



THE HARVARD CHAN
MICROBIOME IN
PUBLIC HEALTH CENTER



Microbial community

(metagenomic)

functional profiling

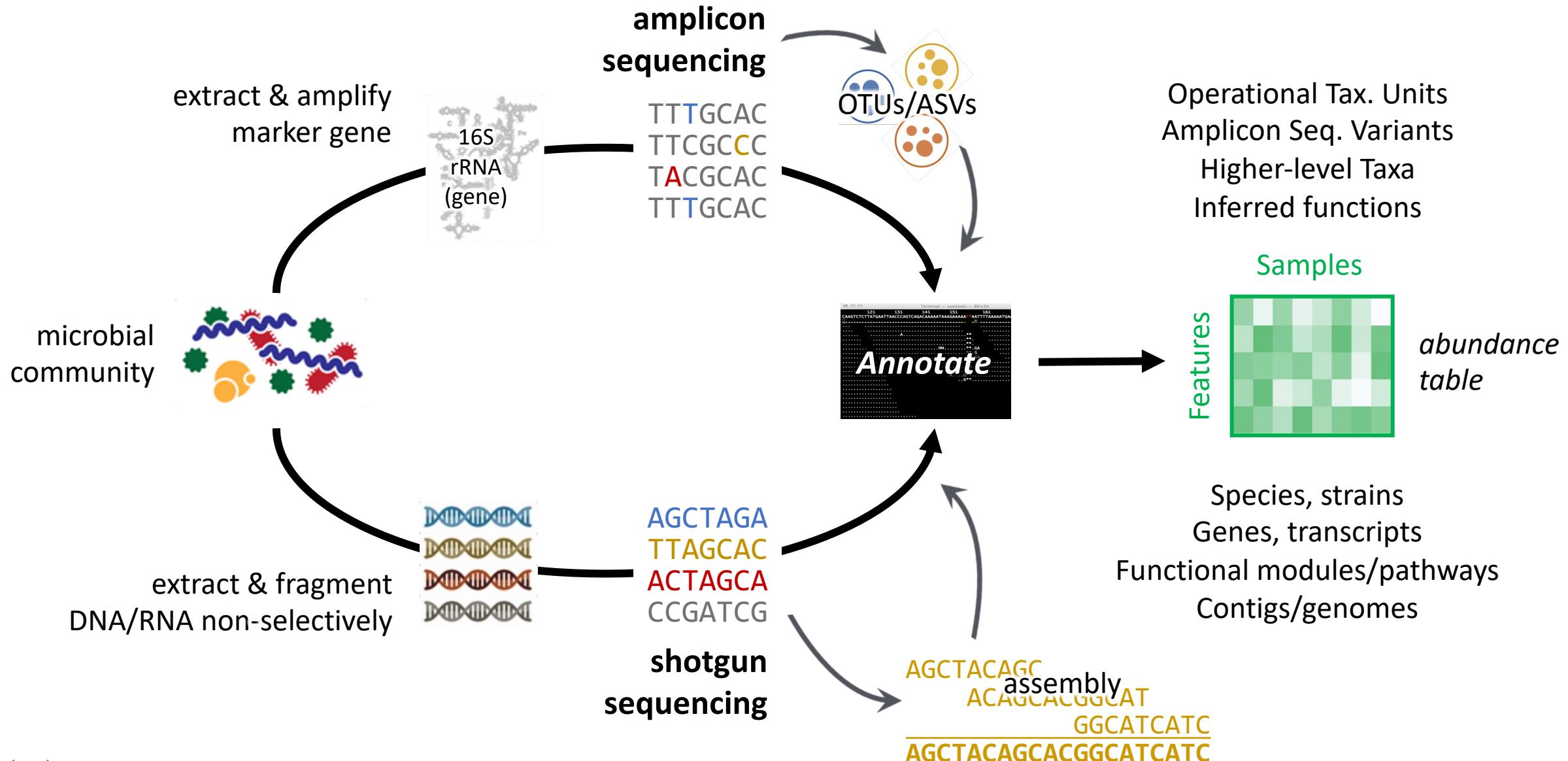


Summary

- Introduction to functional profiling
- Tiered search with HUMAnN
- Core functions of the human microbiome
- Contributional diversity
- What's new in bioBakery 3.0?
- Prioritizing bioactive gene families with MetaWIBELE



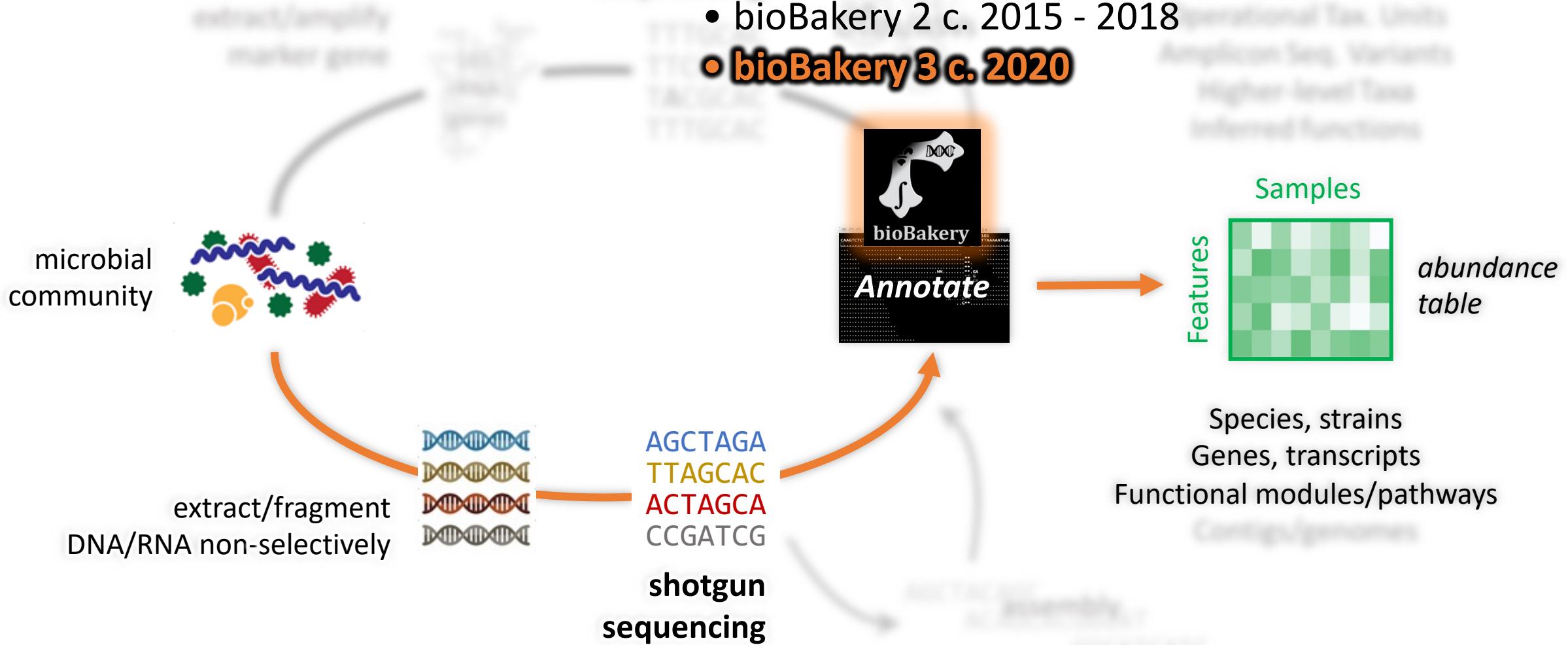
Everything I learned so far at STAMPS in one slide





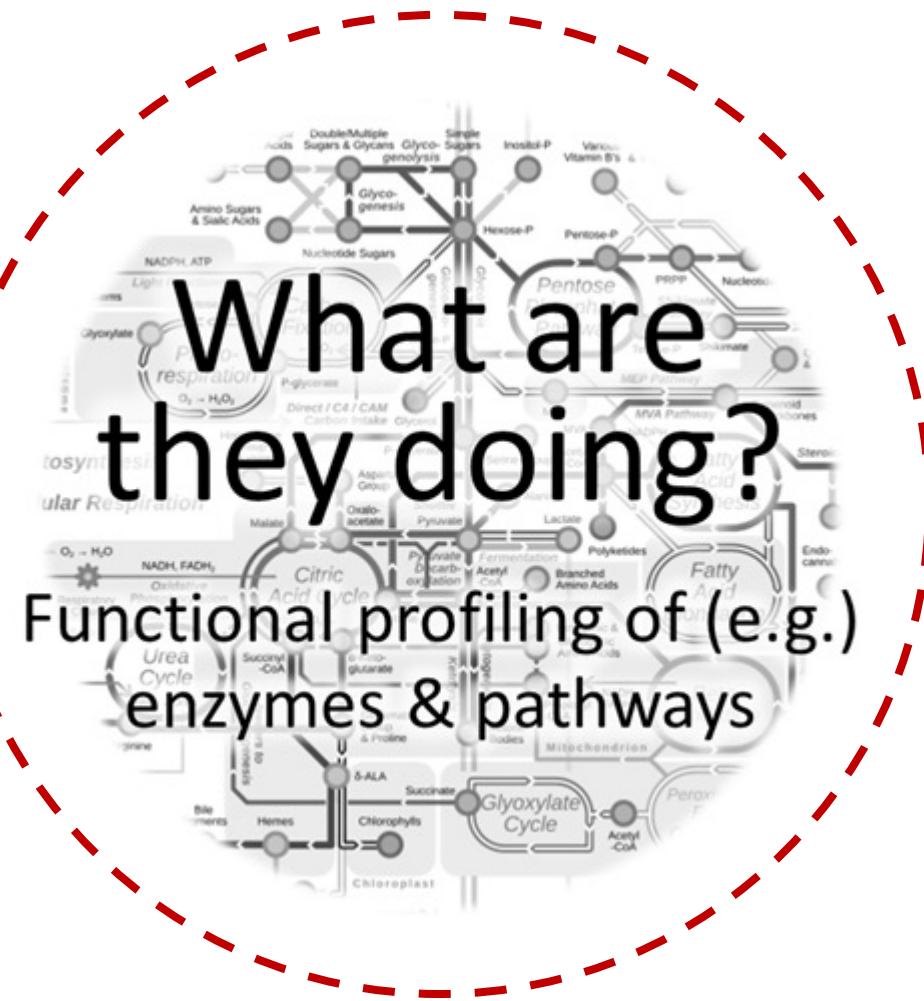
bioBakery 3 directly annotates meta'omic reads

- bioBakery 1 c. 2012
- bioBakery 2 c. 2015 - 2018
- **bioBakery 3 c. 2020**





Two big questions of shotgun meta'omics





What is a functional profile?

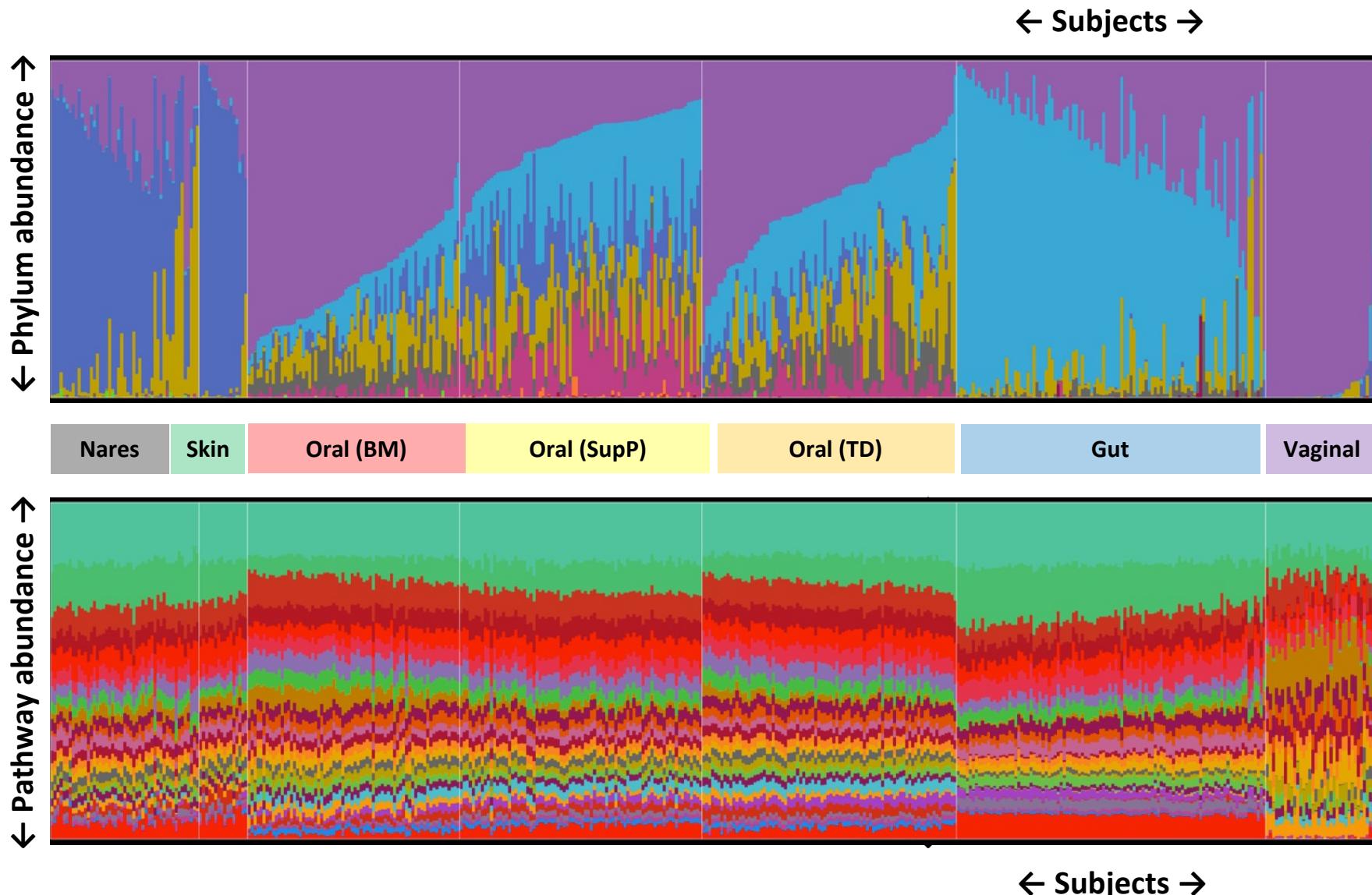


- A functional profile is a vector of relative abundance values describing the molecular functions of a microbial community biosample
- The abundances correspond to features at a fixed level of resolution
 - Millions of conserved gene families (e.g. UniRef90)
 - 10Ks of broad gene families / enzymatic functions (e.g. UniRef50, KO, eggNOG)
 - 100s of metabolic pathways (e.g. KEGG, MetaCyc)
- The abundances form a composition (i.e. they sum to 100%)
- Abundances are proportional to function coverage in the sequencing pool
- Abundances are proportional to copy number in the underlying biosample



Why care about microbiome function?

HMP,
Nature, 2012



Community success
in an environment
requires being able
to do specific tasks
(functions)

Different “bags of
functions” (i.e. taxa)
can collectively
satisfy requirements

Thus, community
function is more
conserved than
structure/taxonomy



Two big questions of shotgun meta'omics

Who is there?

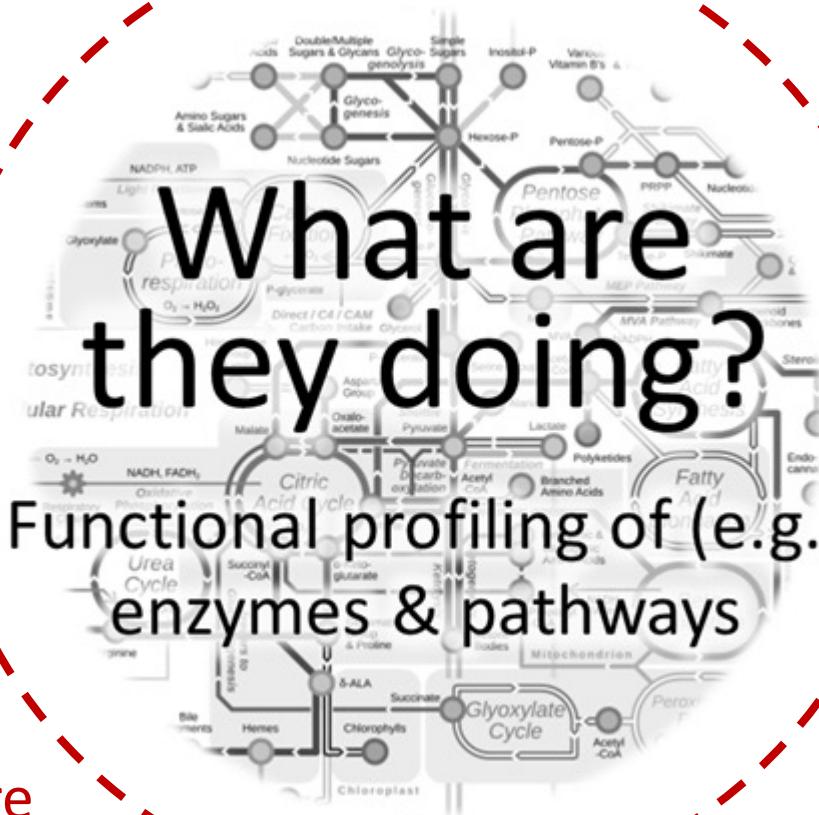
Taxonomic profiling of (e.g.) phyla, species & strains



We often want to answer these at the same time

What are they doing?

Functional profiling of (e.g.) enzymes & pathways



Which organisms are supplying which functions?

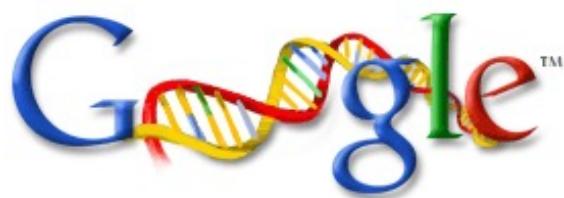


Functional profiling by homology-based search



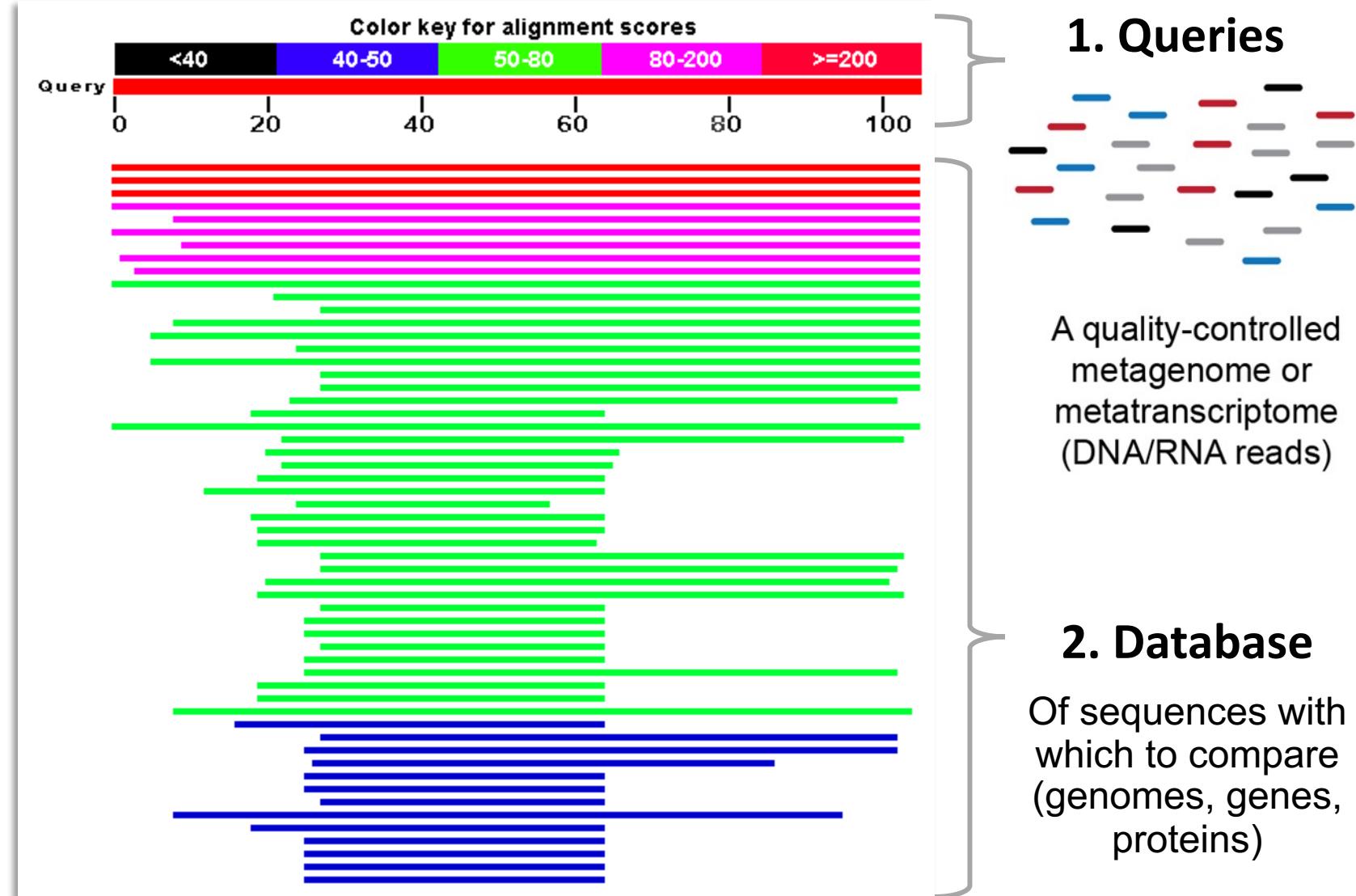
Profiling meta'omes by homology-based search

Three key ingredients:



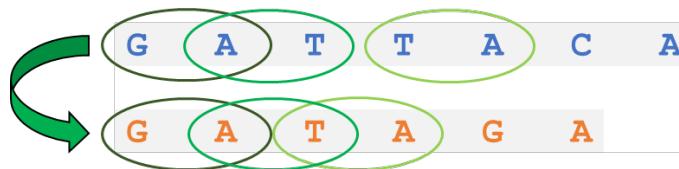
3. Search “Engine”

Of which BLAST is the prototypical example, though we often prefer other, more specialized methods





3. Search engines: Alignment and mapping



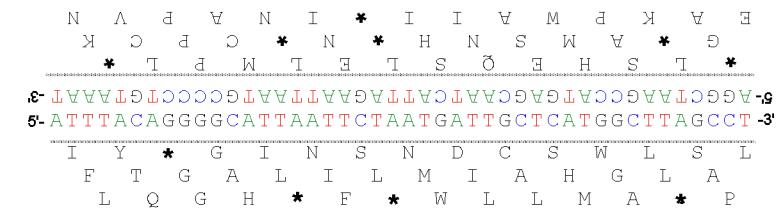
A *read mapping* places reads into source sequences without using alignment (usually based on overlap of k -mer words)

Example Tools

Kallisto, Salmon (general mappers)
Clark, Kraken (taxonomic binners)



A *nucleotide alignment* maps homologous sites between the query and database sequences (inserting gaps to help)



A *translated alignment* compares the 6 possible translations of a DNA sequence to a protein sequence database

Faster

More specific (fewer false positives)

Less sensitive (more false negatives)



Slower

Less specific (more false positives)

More sensitive (fewer false negatives)



3. Search engines: Alignment and mapping



A *read mapping* places reads into source sequences without using alignment (usually based on a library of known words)

However, this sensitivity makes translated search more prone to false positive hits

Kallisto, Salmon (general mappers)

Clark, Kraken (taxonomic binners)

And, in addition, it is very slow compared to nucleotide-level mapping/alignment
(6 protein searches per nucleotide read)

Less sensitive (more false negatives)



A nucleotide alignment
maps homologous sites between the query and database sequences
(Inserting gaps to help)

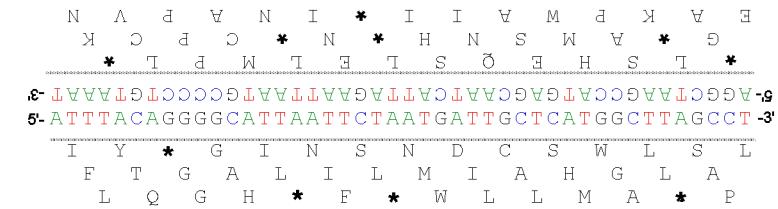
blastn (for fewer, longer sequences)

bowtie2, bwa-mem (for short reads)

Faster

More specific (fewer false positives)

Less sensitive (more false negatives)



A translated alignment compares the 6 possible translations of a DNA sequence to a protein sequence database

Example Tools

blastp (classic but slow)

DIAMOND, MMSeqs2 (fast)

Slower

Less specific (more false positives)

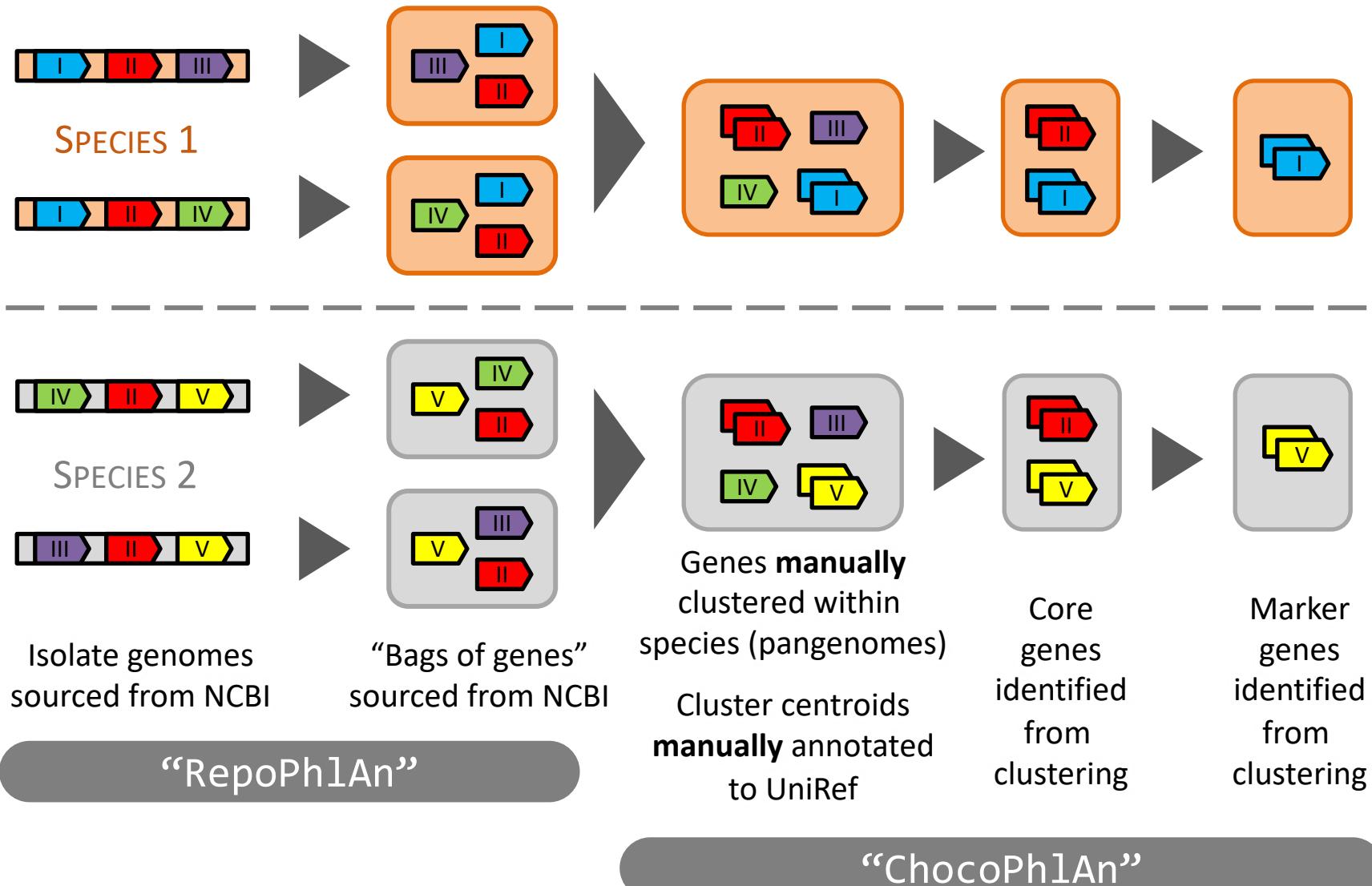
More sensitive (fewer false negatives)



The HUMAnN approach to functional profiling



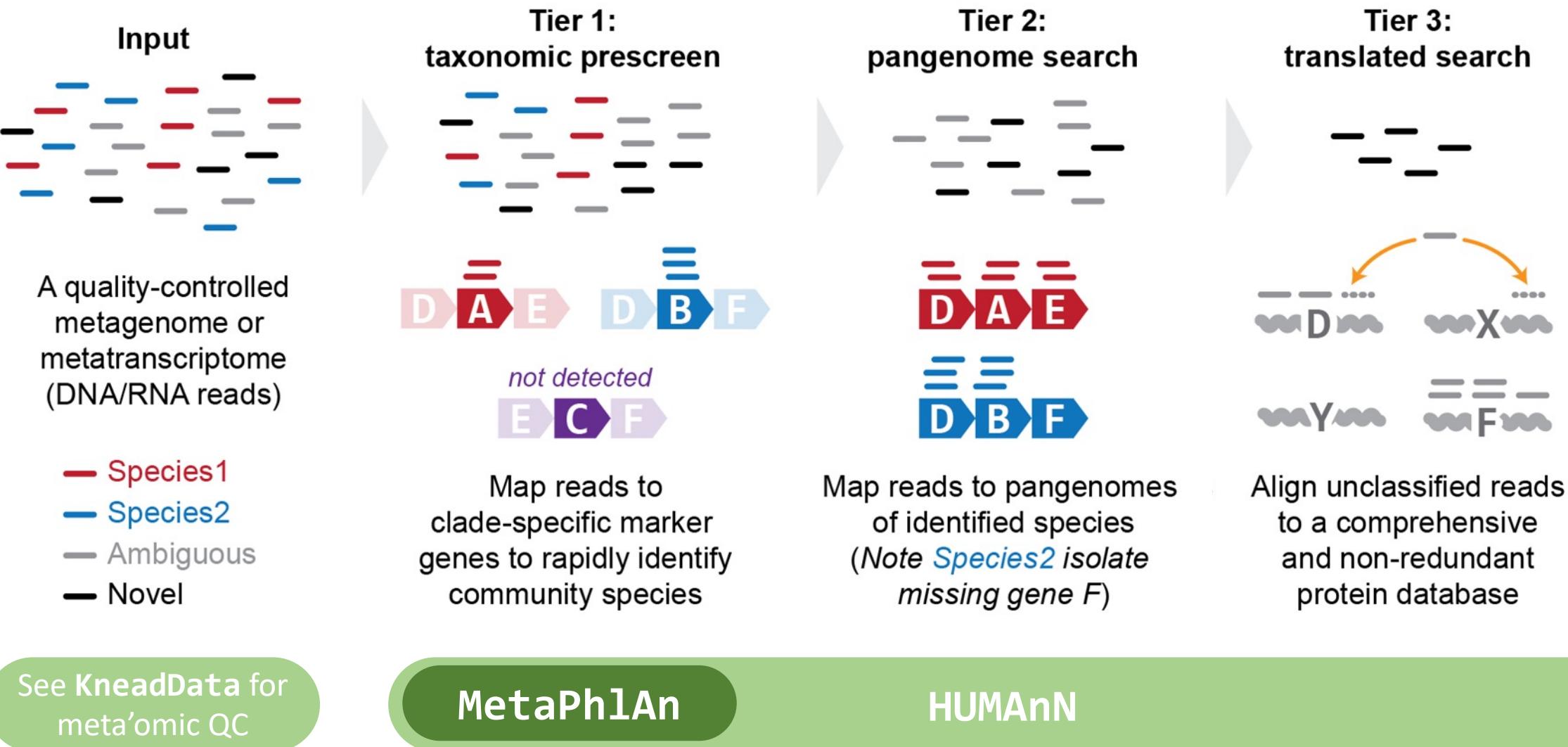
HUMAnN takes advantage of functionally annotated pangenomes and marker-gene based species detection



- RepoPhlAn and ChocoPhlAn are database systems for identifying markers genes from isolate genomes (among other tasks).
- We will use other products of these systems in subsequent analysis methods.

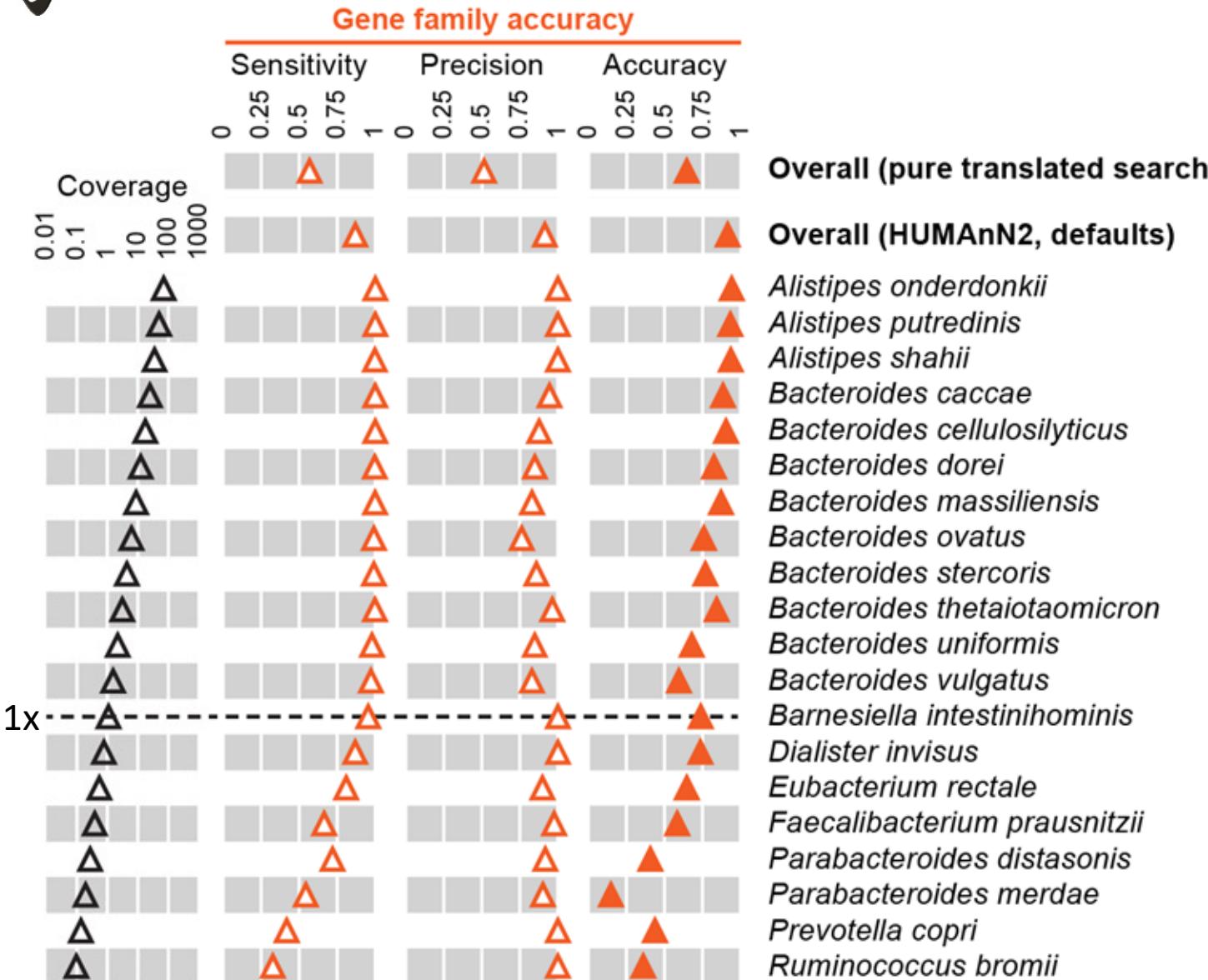


Accurate, efficient meta'omic profiling by tiered search





HUMAnN's tiered search is a more accurate approach



- Evaluation on a synthetic gut metagenome (top-20 HMP)
- 0.1x to 100x coverage
- Compare expected and observed gene coverage
- Pure translated search loses specificity (false positives)
- HUMAnN's tiered search is more accurate overall
- Per-species results for free!
 - Specific for all species
 - sensitive when coverage > 1x

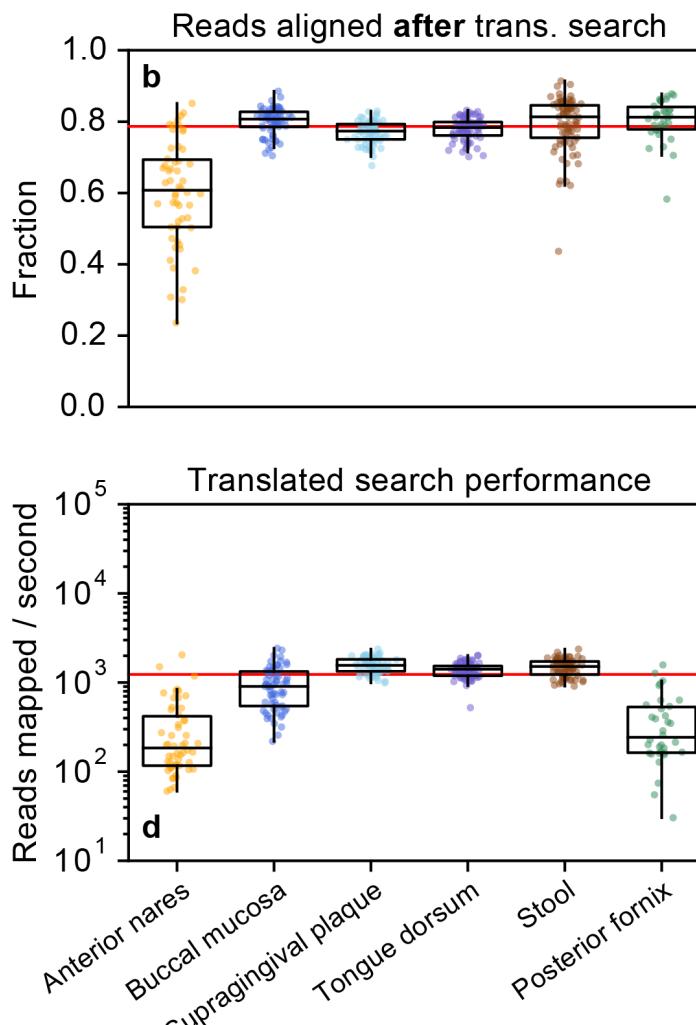
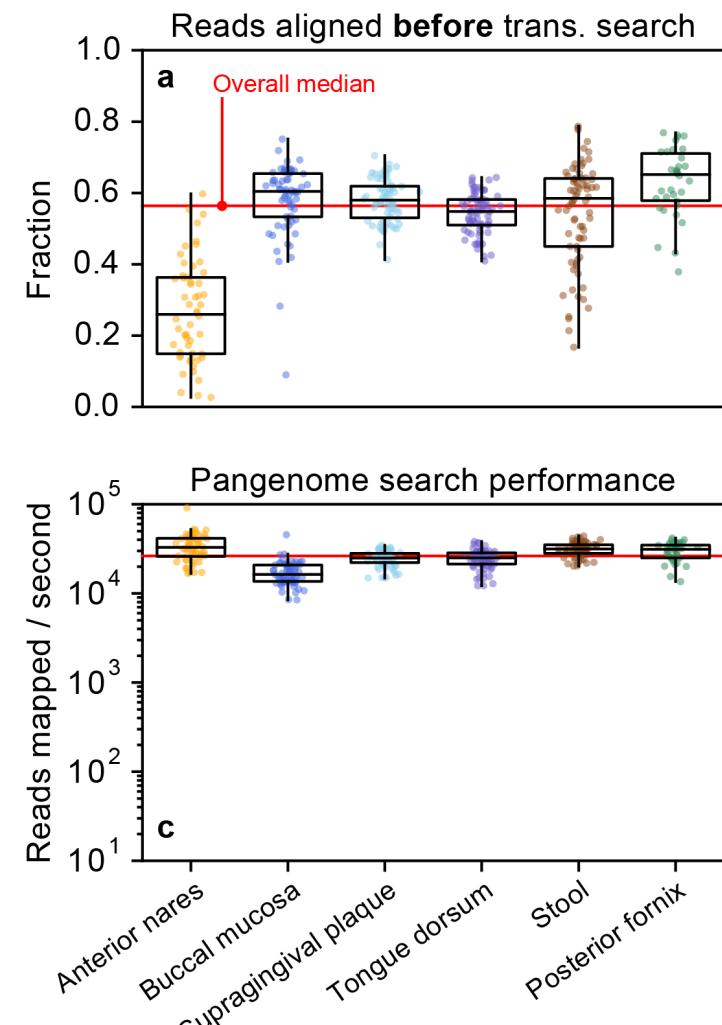


Tiered search is also faster!

Lots of reads explained during fast search (or this would be a moot point)



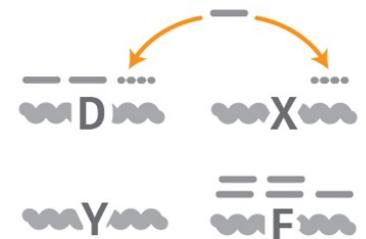
Reads align 10x faster to selected pangenomes vs. translated search



If $1/T$ reads unaligned after pangenome search, then HUMAnN is T times faster than pure translated search (e.g. $\frac{1}{2} \rightarrow 2x$ faster)

Additional reads align in translated search, especially those from skin

Consequence of a smaller database and nucleotide-level alignment





HUMAnN's species-stratified abundance format

Tiered search is *faster, more accurate, and also provides more biological information*

UniRef90_R6K3Z5: IMP dehydrogenase	600.95
UniRef90_R6K3Z5: IMP dehydrogenase Bacteroides_caccae	234.76
UniRef90_R6K3Z5: IMP dehydrogenase Bacteroides_dorei	107.38
UniRef90_R6K3Z5: IMP dehydrogenase Bacteroides_ovatus	92.18
UniRef90_R6K3Z5: IMP dehydrogenase Bacteroides_stercoris	83.95
UniRef90_R6K3Z5: IMP dehydrogenase Bacteroides_vulgatus	57.27
UniRef90_R6K3Z5: IMP dehydrogenase unclassified	25.41



HUMAnN's species-stratified abundance format

UniRef gene family ID	UniRef gene family name	Community-total gene abundance (RPK)
UniRef90_R6K3Z5: IMP dehydrogenase		600.95
UniRef90_R6K3Z5: IMP dehydrogenase <i>Bacteroides caccae</i>		234.76
UniRef90_R6K3Z5: IMP dehydrogenase <i>Bacteroides dorei</i>		107.38
UniRef90_R6K3Z5: IMP dehydrogenase <i>Bacteroides ovatus</i>		92.18
UniRef90_R6K3Z5: IMP dehydrogenase <i>Bacteroides stercoris</i>		83.95
UniRef90_R6K3Z5: IMP dehydrogenase <i>Bacteroides vulgatus</i>		57.27
UniRef90_R6K3Z5: IMP dehydrogenase unclassified		25.41

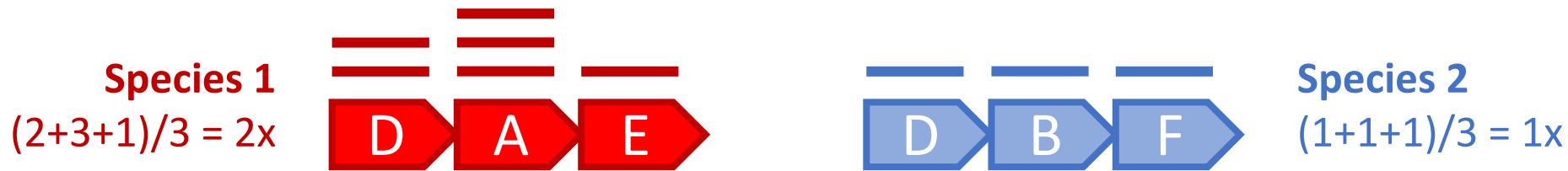


Per-species & unclassified stratifications



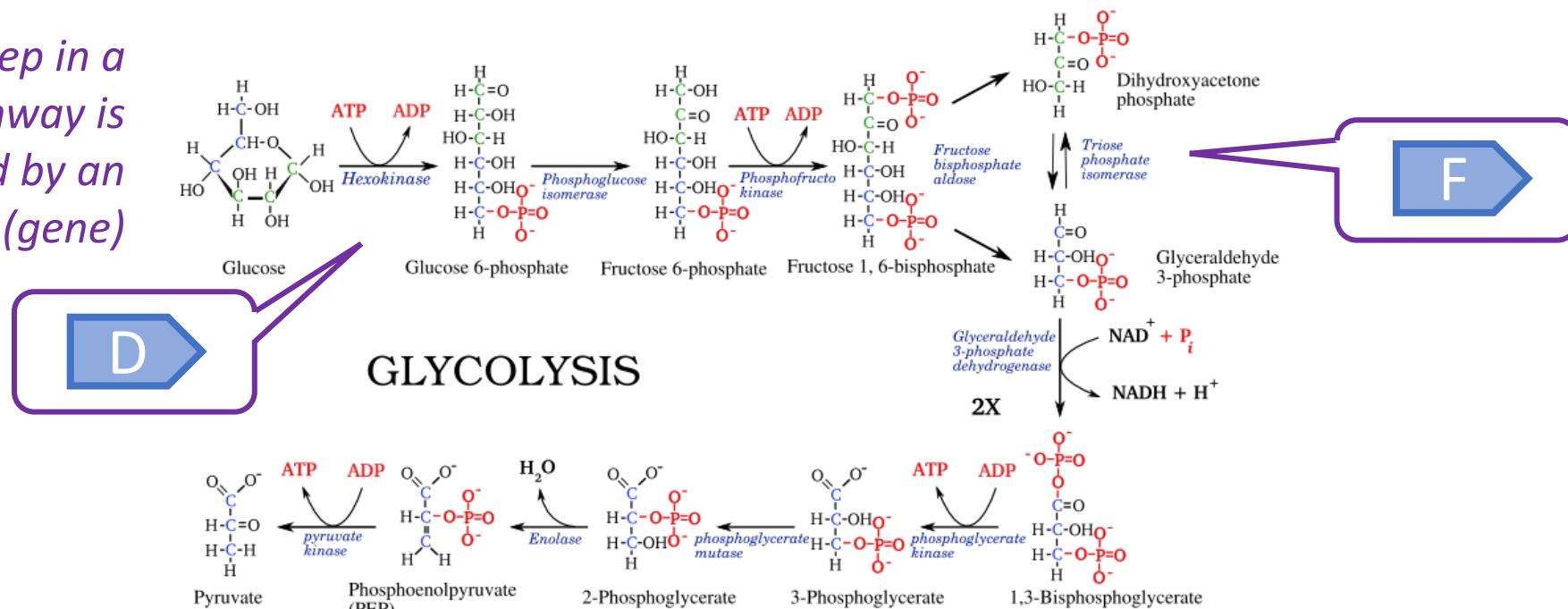
Regrouping genes to pathways/modules

- If a species is present, we expect to see its marker genes evenly covered:



- If a pathway is present, we expect to see its enzymes evenly covered:

Each step in a metabolic pathway is governed by an enzyme (gene)

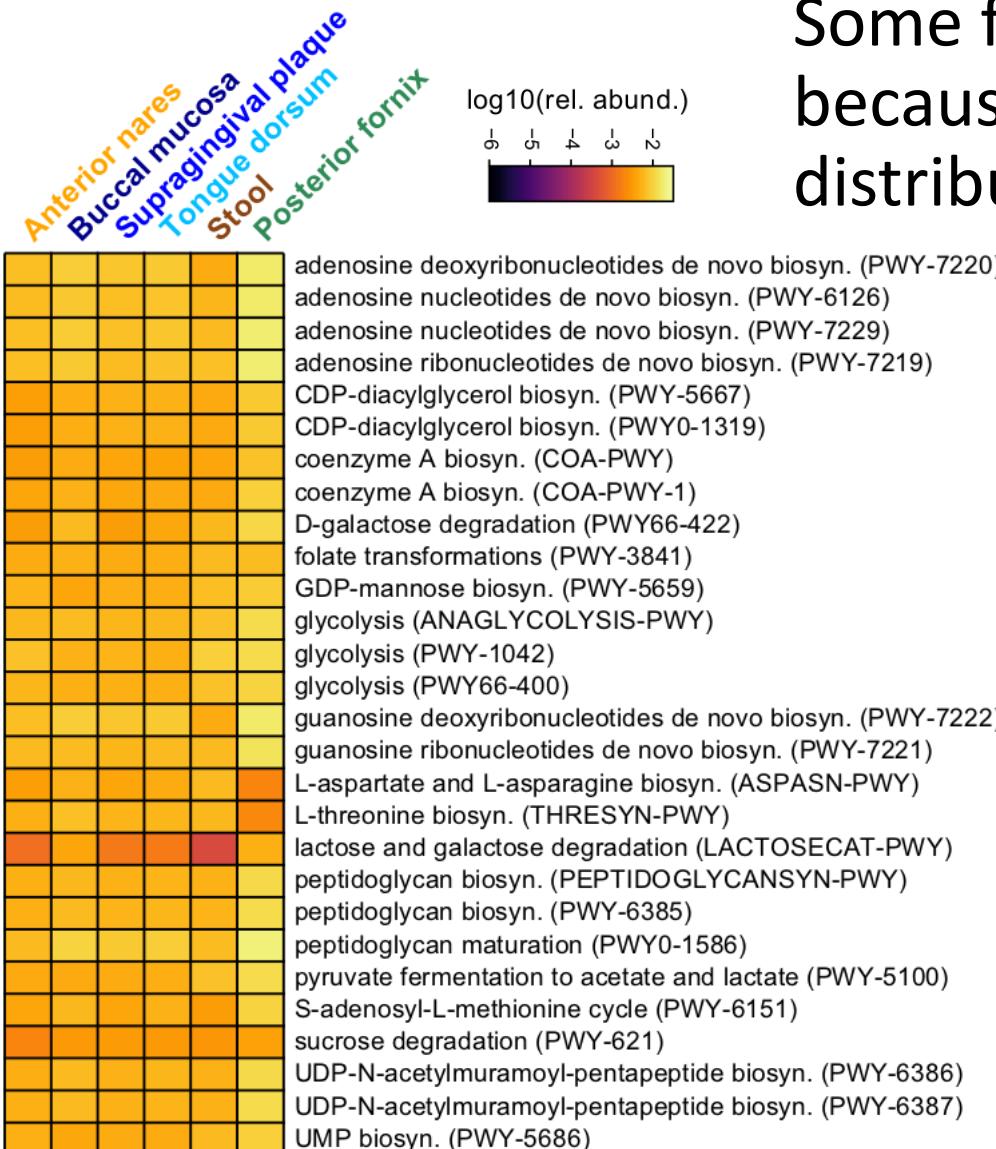




Core functions of the
human microbiome &
contributional diversity



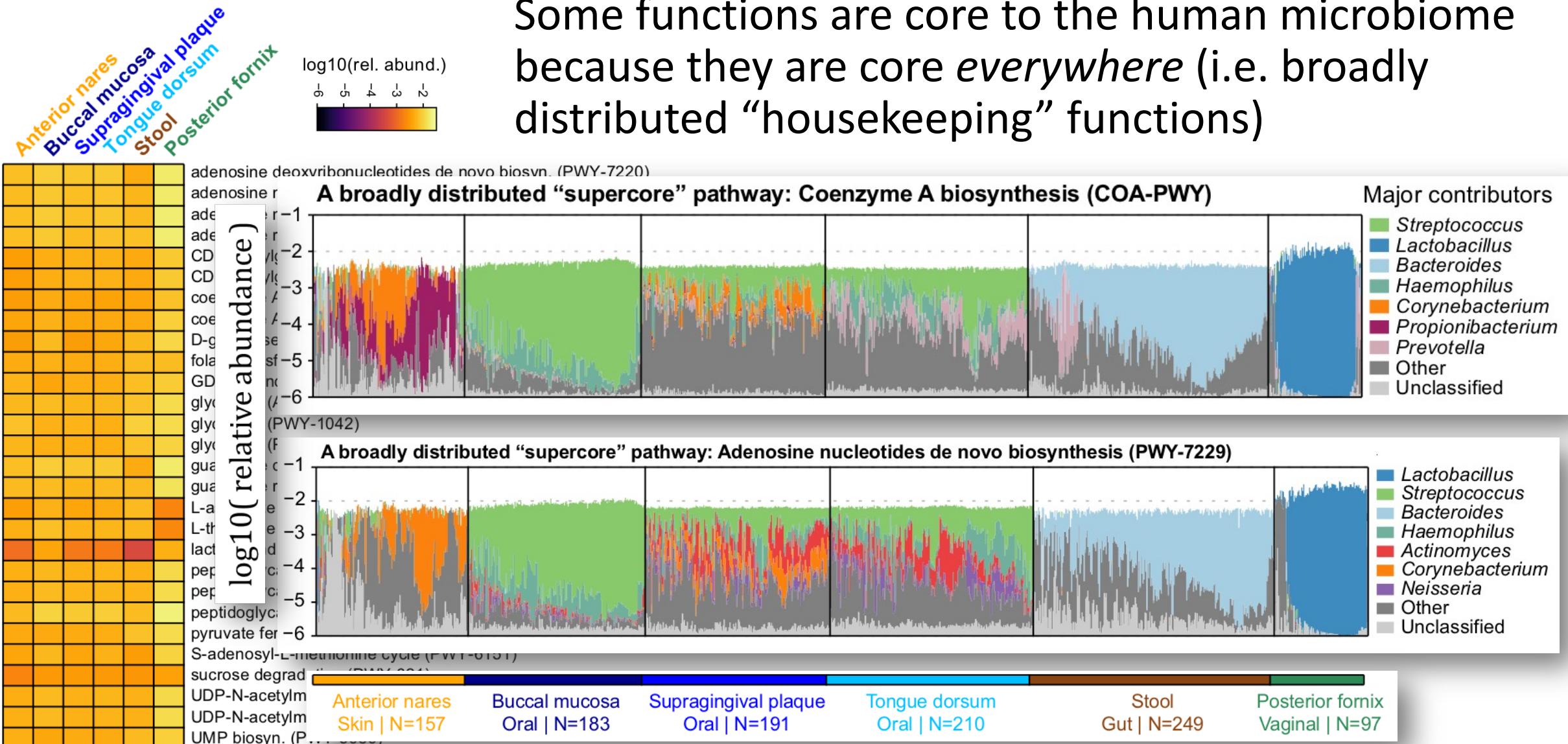
Core functions of the human microbiome



Some functions are core to the human microbiome because they are core *everywhere* (i.e. broadly distributed “housekeeping” functions)

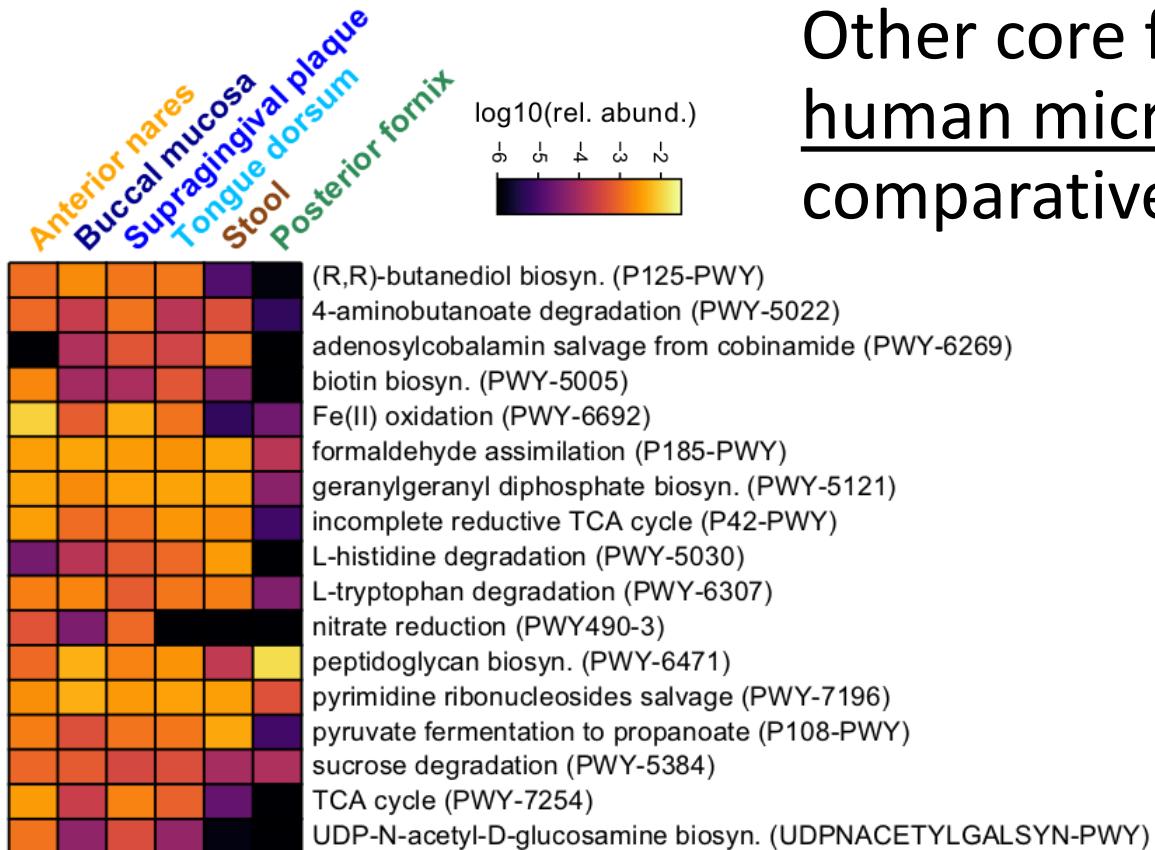


Core functions of the human microbiome

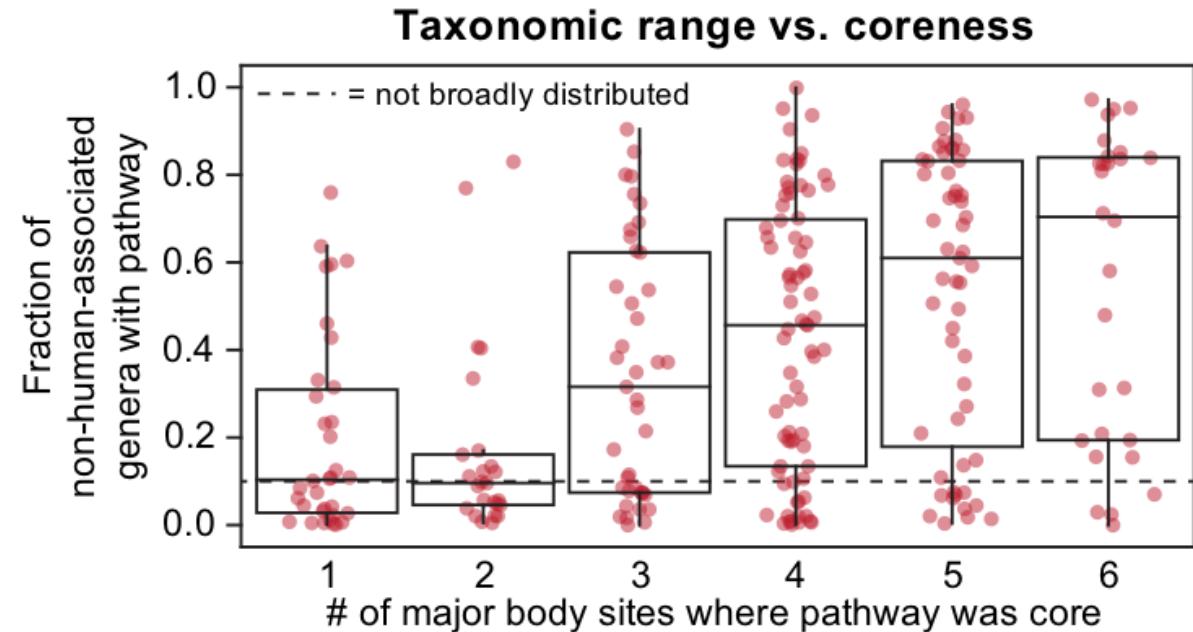




Core functions of the human microbiome



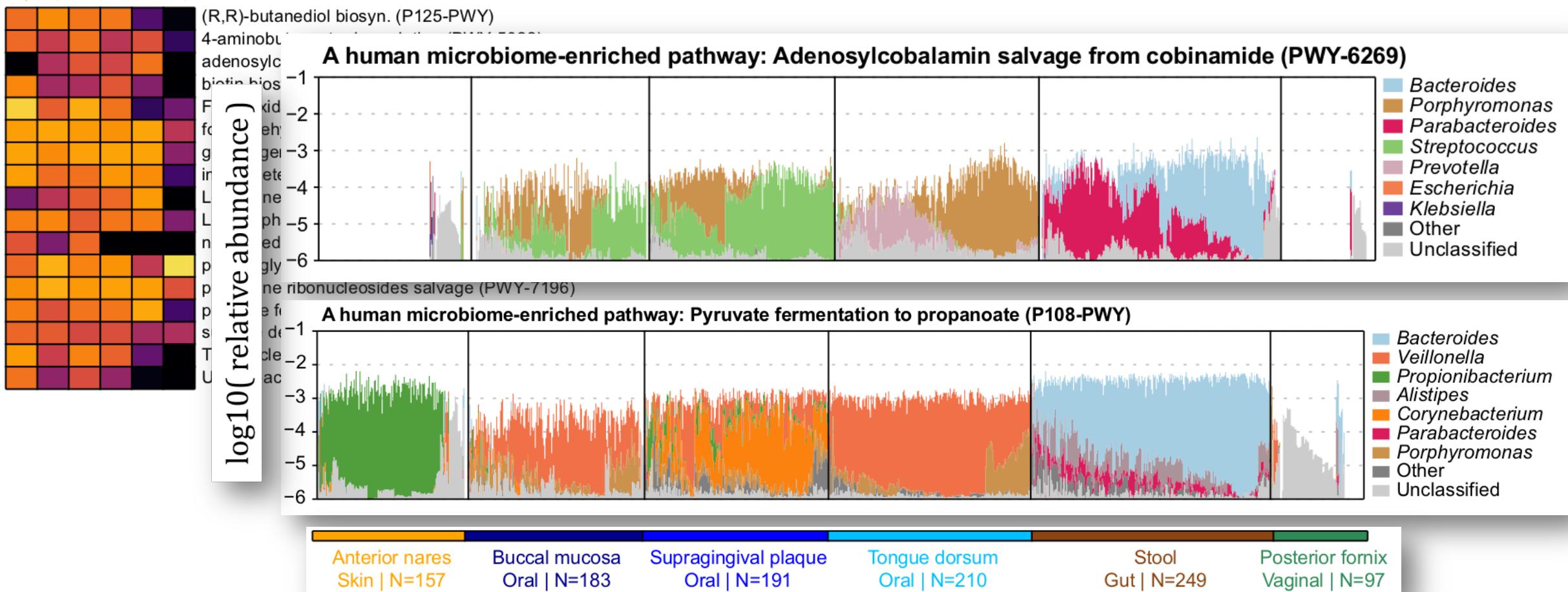
Other core functions are specifically enriched in the human microbiome (core at multiple body sites and comparatively rare in non-human-associated bugs)



Core functions of the human microbiome

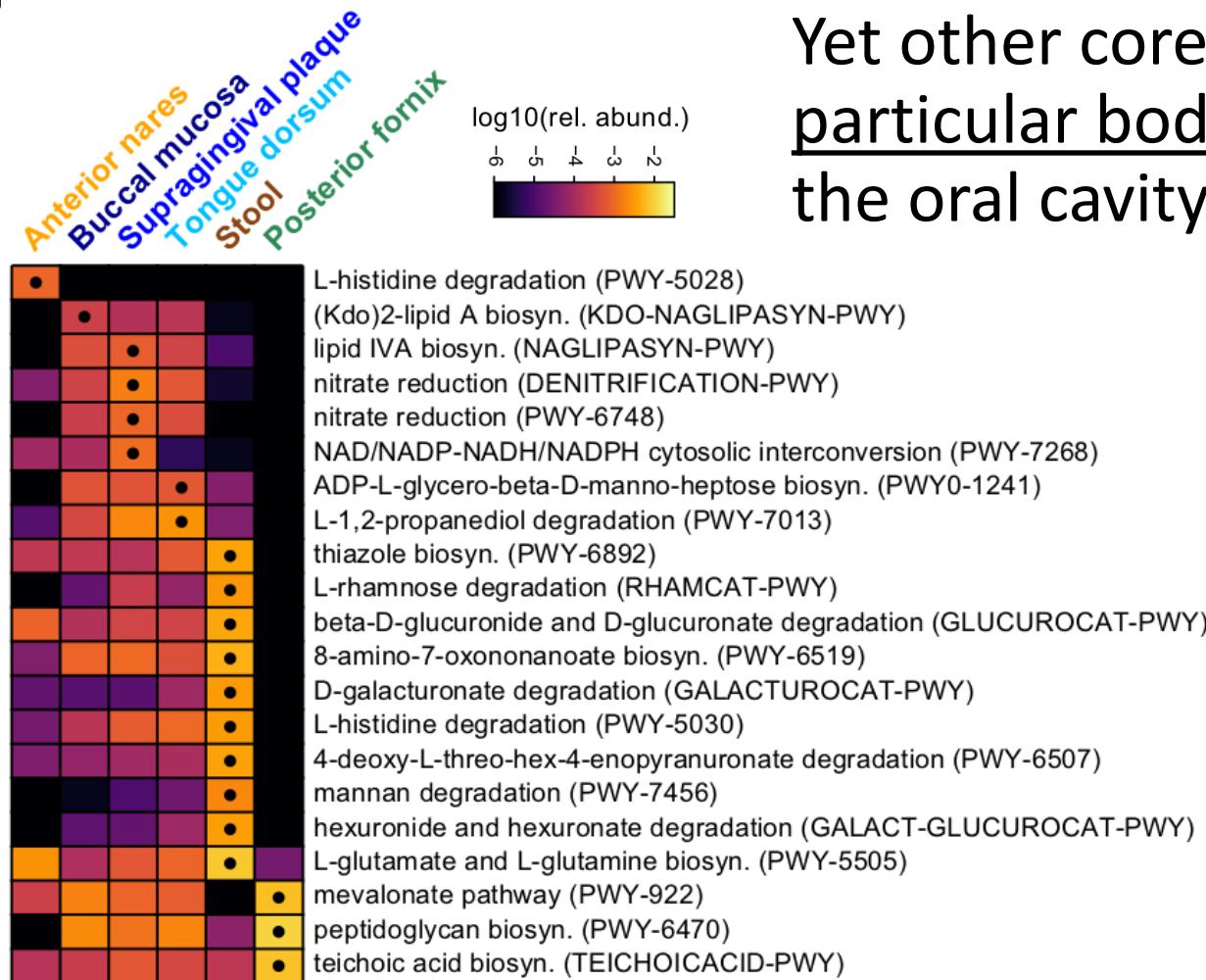


Other core functions are specifically enriched in the human microbiome (core at multiple body sites and comparatively rare in non-human-associated bugs)





Core functions of the human microbiome

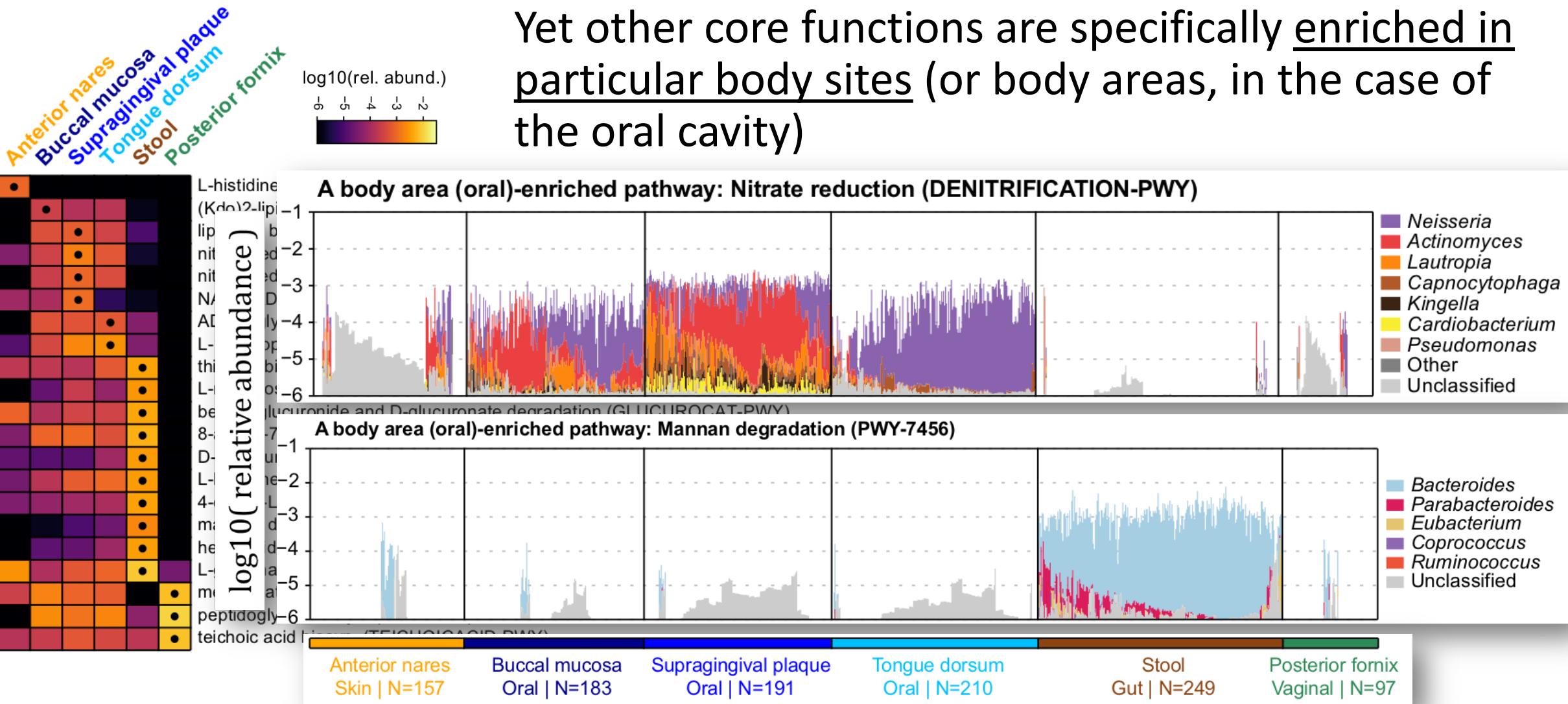


Yet other core functions are specifically enriched in particular body sites (or body areas, in the case of the oral cavity)



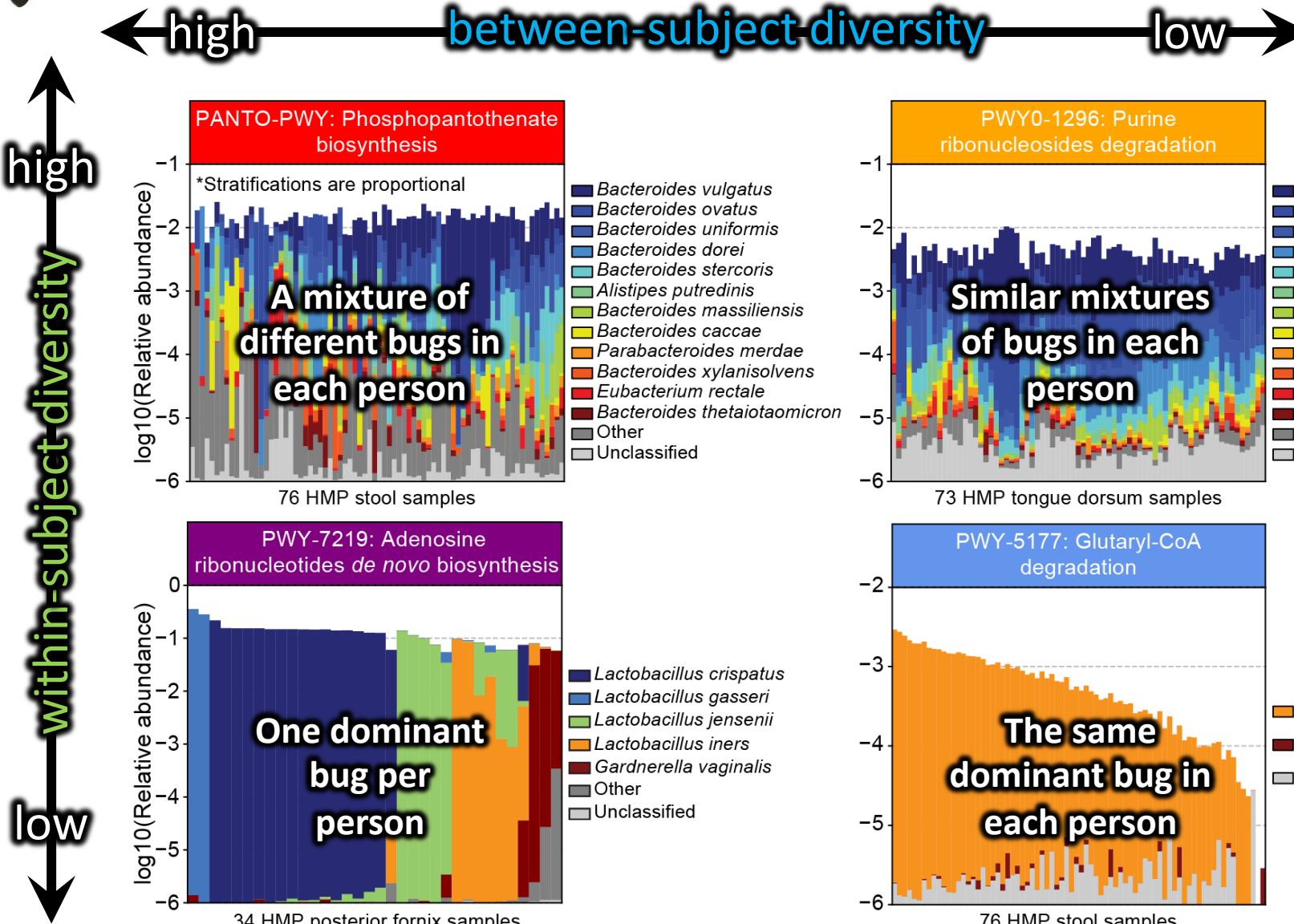
Core functions of the human microbiome

Lloyd-Price,
Nature, 2017





Contributional diversity of core microbiome pathways



Apply community level diversity measures to species' contributes to individual functions

contributional diversity

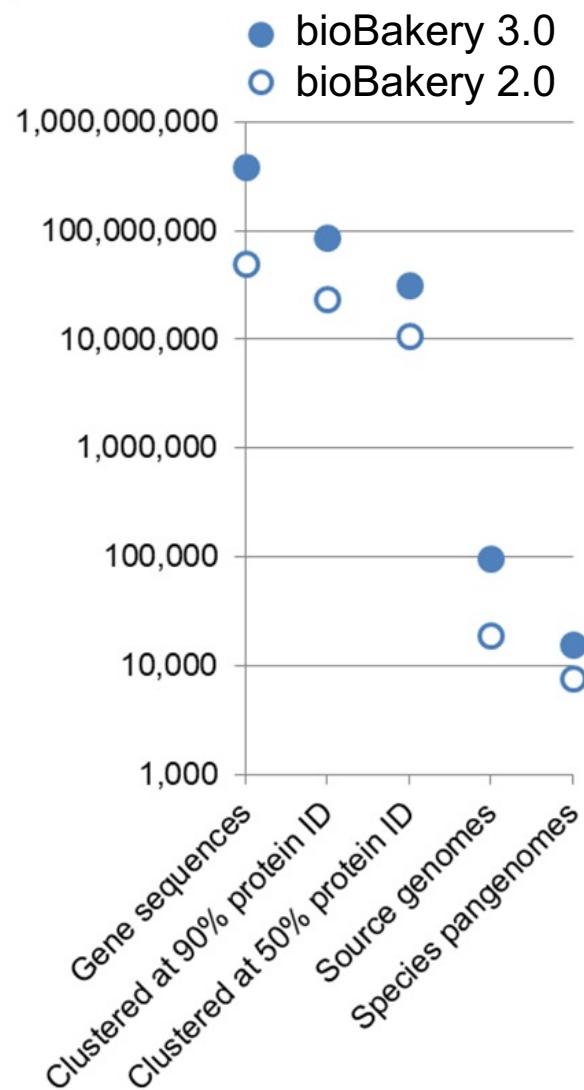
Within (alpha) and between (beta) sample flavors



What's new in bioBakery 3.0? (including HUMAnN 3.0)



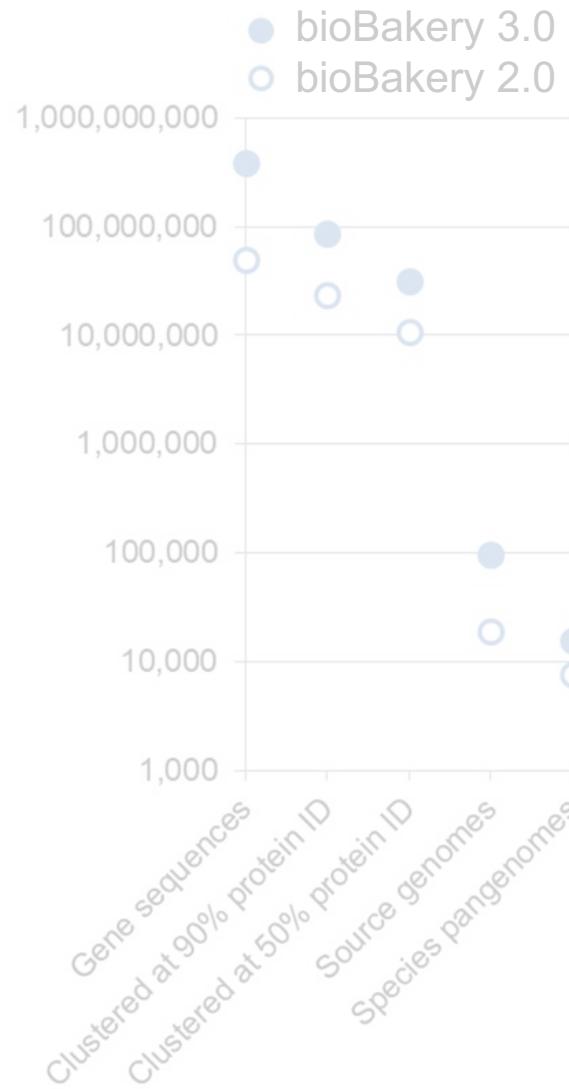
bioBakery 3 builds bigger databases more easily



- 99K microbial genomes (**5x increase** vs. bioBakery 2)
- 16K microbial species (**3x increase**)
- 87M gene families (**3x increase**)
- 389M gene sequences (**8x increase**)

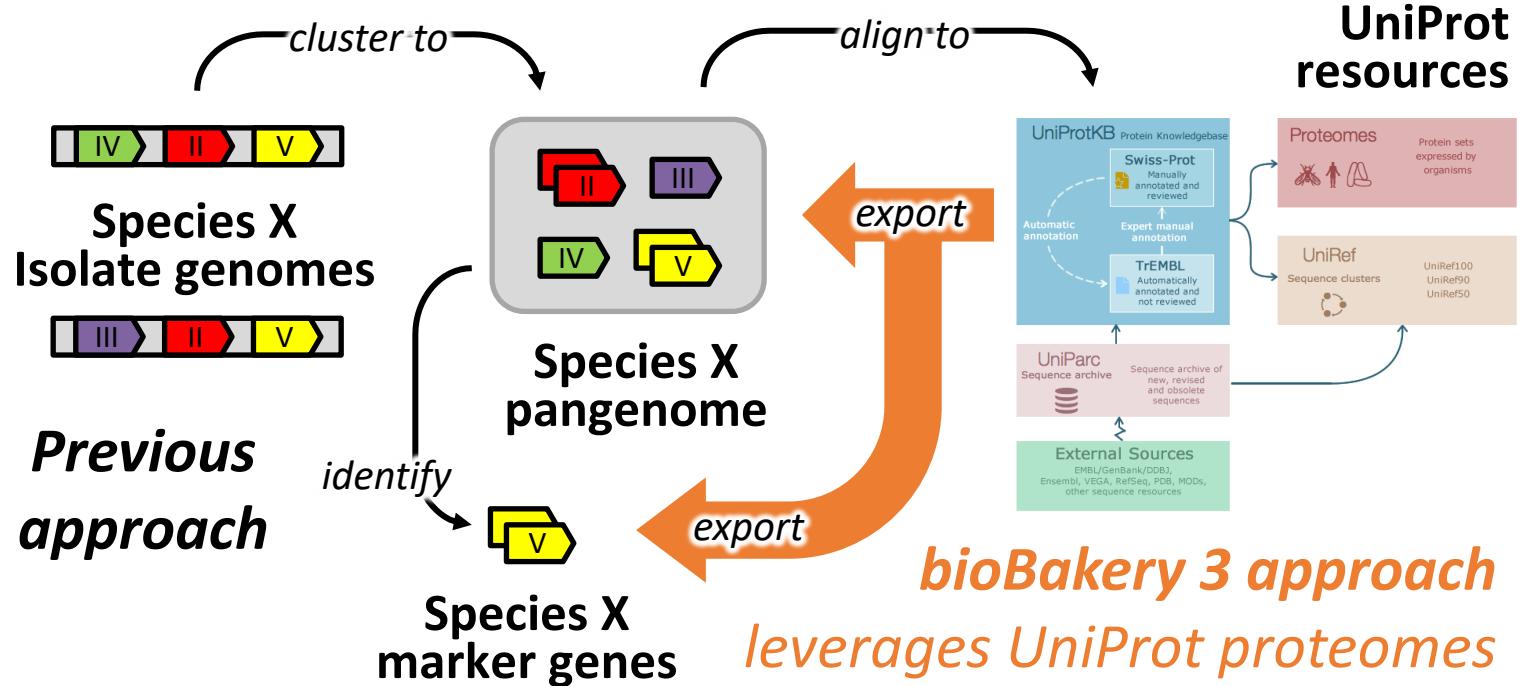


bioBakery 3 builds bigger databases more easily



- 99K microbial genomes (**5x increase vs. bioBakery 2**)
- 16K microbial species (**3x increase**)
- 87M gene families (**3x increase**)
- 389M gene sequences (**8x increase**)

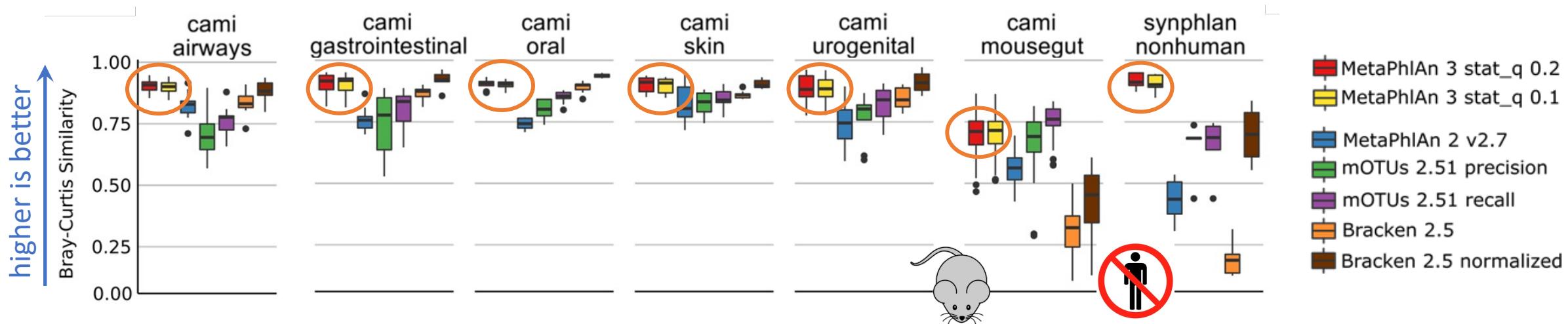
Facilitated by new methods for sequence database construction



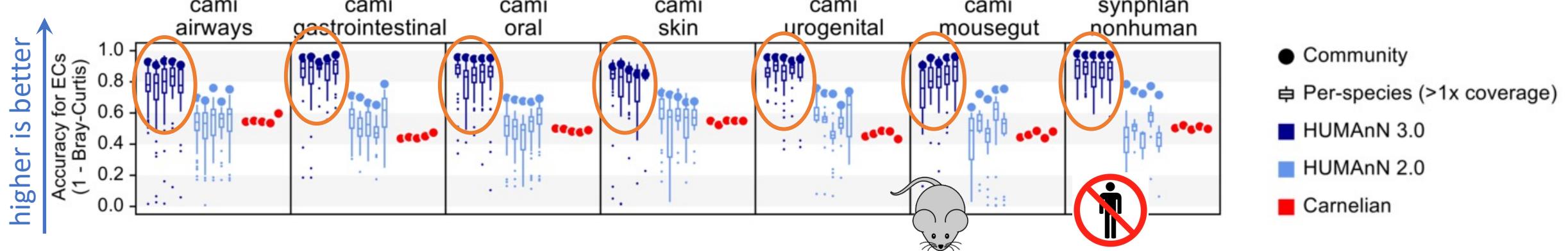


bioBakery 3 accurately profiles synthetic communities

- MetaPhlAn 3 accurately profiles species from synthetic metagenomes (CAMI, etc.)



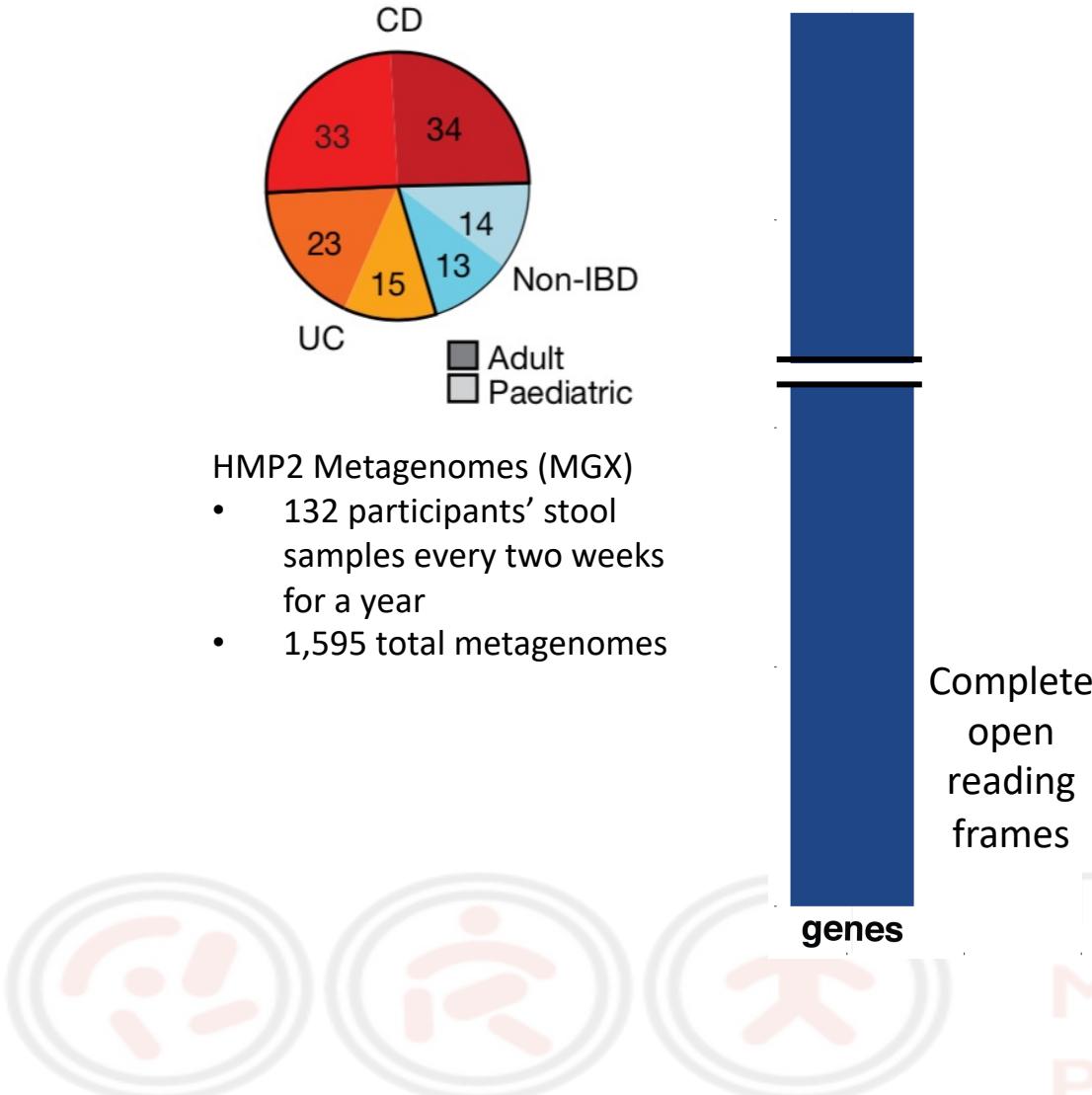
- HUMAnN 3 accurately profiles enzyme abundance (total and per-species)





Prioritizing bioactive gene families with MetaWIBELE

Even in the human gut, uncharacterized proteins abound



44,585,131



Complete
open
reading
frames

2,088,122

gene catalogs

95%
nucleotide
identity
90% coverage

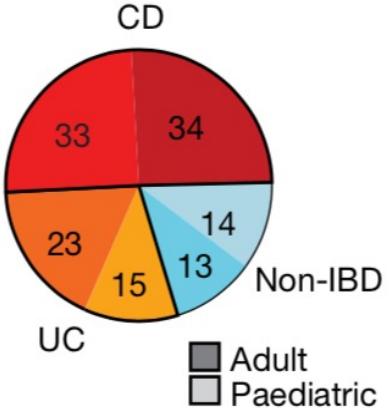
1,665,233

protein families

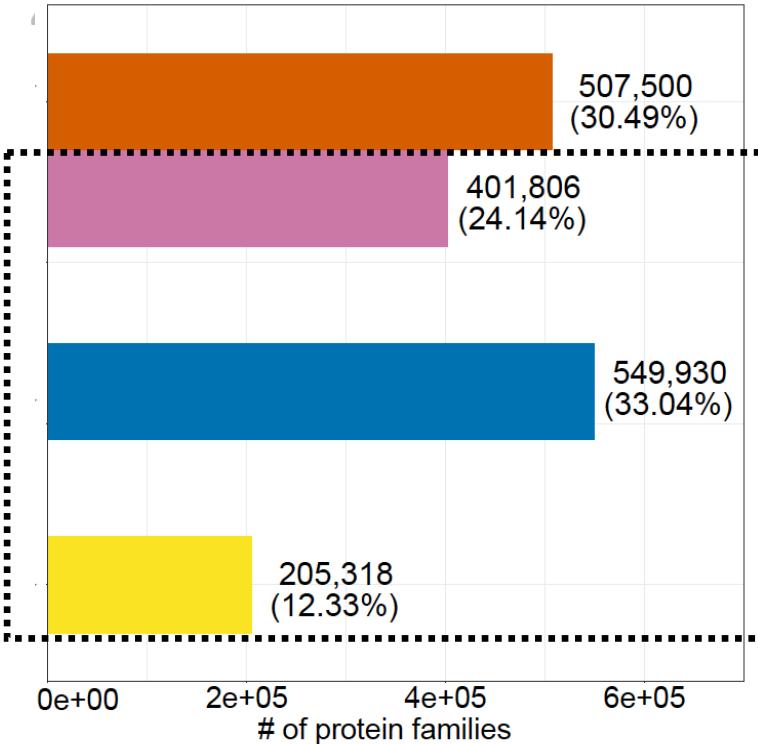
90%
amino acid
identity
80% coverage

THE HARVARD
MICROBIOME IN
PUBLIC HEALTH CENTER

Even in the human gut, uncharacterized proteins abound

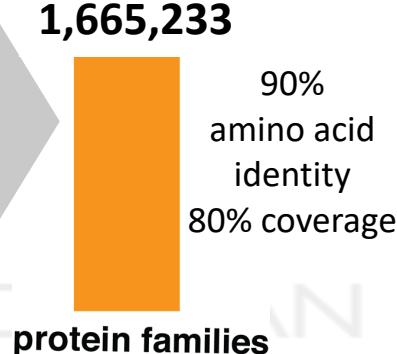


- HMP2 Metagenomes (MGX)
- 132 participants' stool samples every two weeks for a year
 - 1,595 total metagenomes



~70% of protein families
uncharacterized

- Strong homology to known characterized proteins
- Strong homology to known uncharacterized proteins
- Remote homology to known proteins
- No homology to known proteins

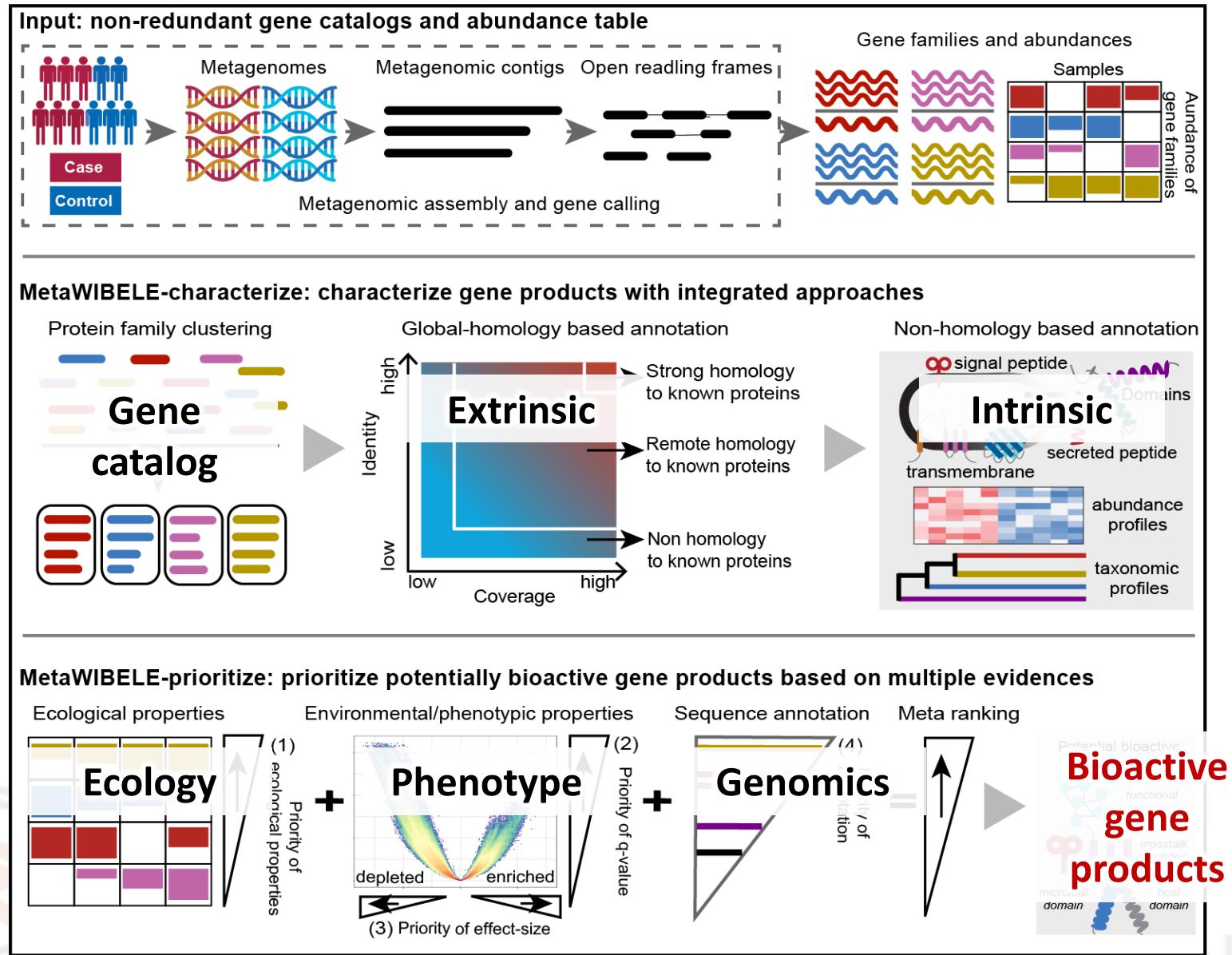


Prioritizing and (initially) characterizing “important” microbial community proteins

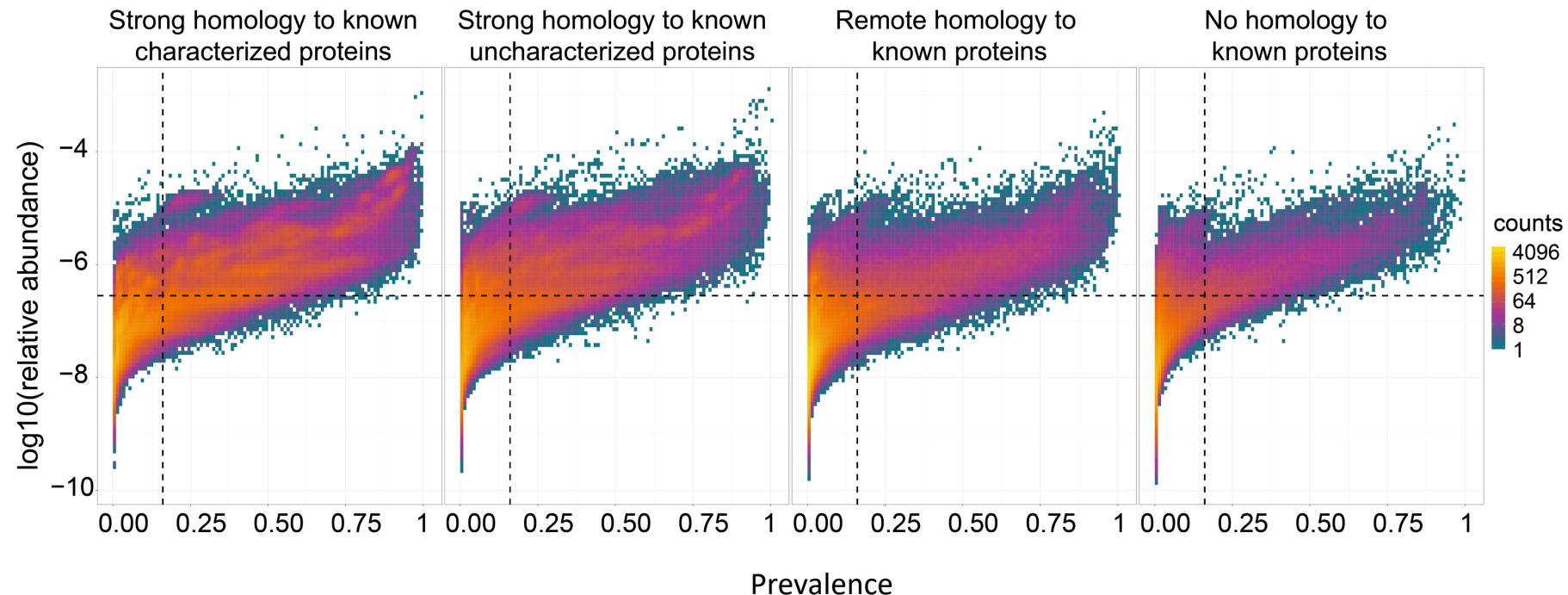


Yancong
Zhang

MetaWIBELE:
Workflow to
Identify novel
Bioactive
Elements in the
microbiome



Uncharacterized proteins are “typical”



THE HARVARD CHAN
MICROBIOME IN
PUBLIC HEALTH CENTER

Bioactivity prioritizations in IBD are supported by external validation

Chaperone-usher pilins (mostly Type I and P pilins) were prioritized during IBD dysbiosis

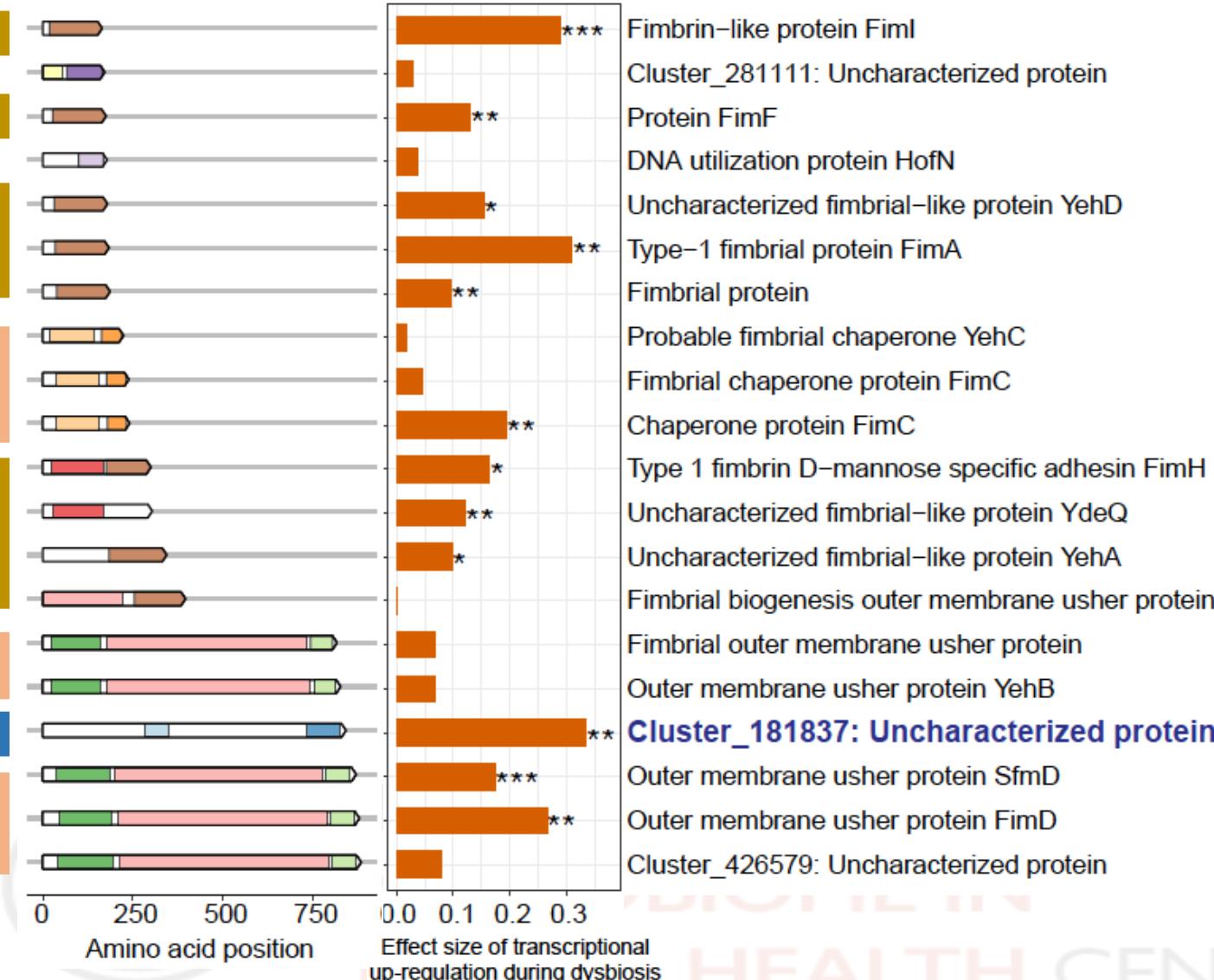
Type I pilin

P pilin

CS1 pilin

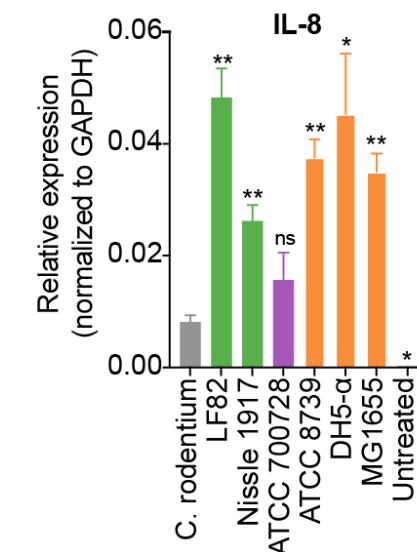
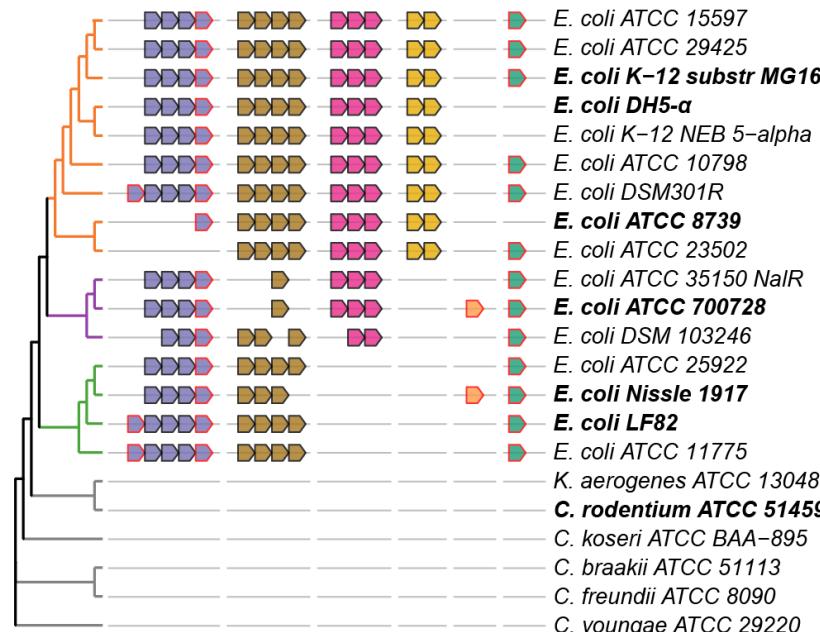
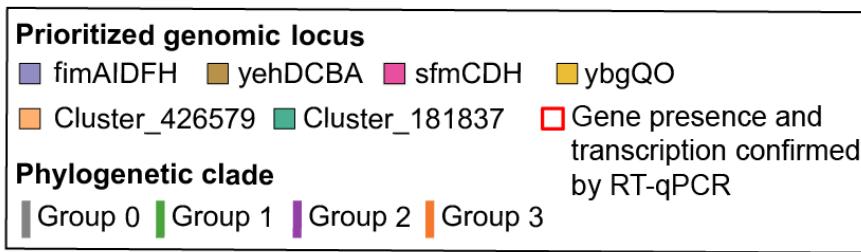
PF00345:Pili and flagellar-assembly chaperone, PapD N-terminal domain
PF00419:Fimbrial protein
PF00577:Outer membrane usher protein
PF01245:Ribosomal protein L19
PF02753:Pili assembly chaperone PapD, C-terminal domain
PF05137:Fimbrial assembly protein (PilN)
PF09160:FimH, mannose binding
PF10610:Thin aggregative fimbriae synthesis protein
PF13953:PapC C-terminal domain
PF13954:PapC N-terminal domain
PF15976:CS1-pili formation C-terminal
PF16967:E-set like domain

* FDR q < 0.05 ** FDR q < 0.01 *** FDR q < 0.001

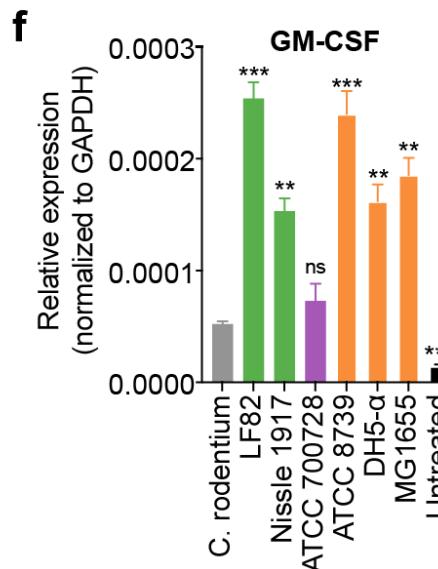


Bioactivity prioritizations in IBD are supported by external validation

Chaperone-usher pilins
(mostly Type I and P
pilins) were prioritized
during IBD dysbiosis

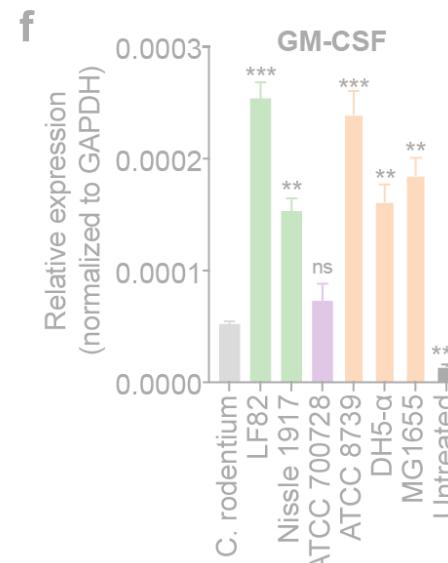
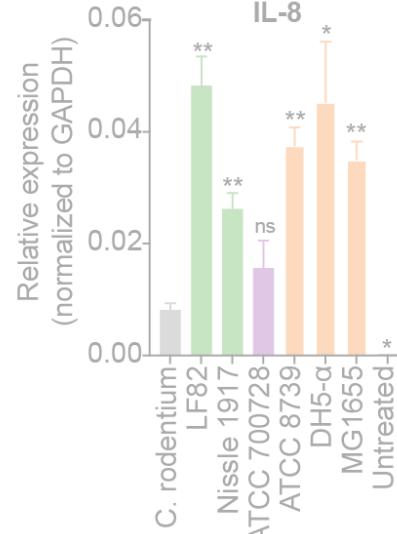
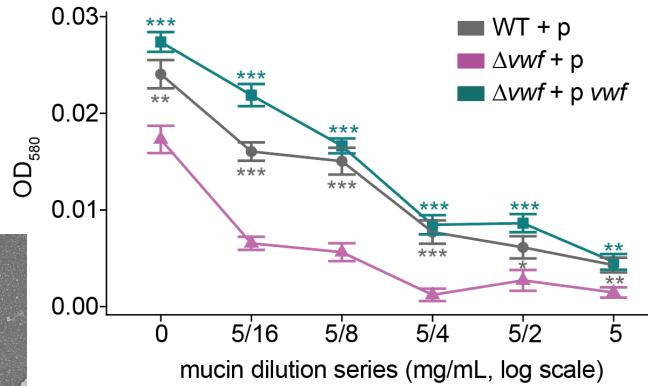
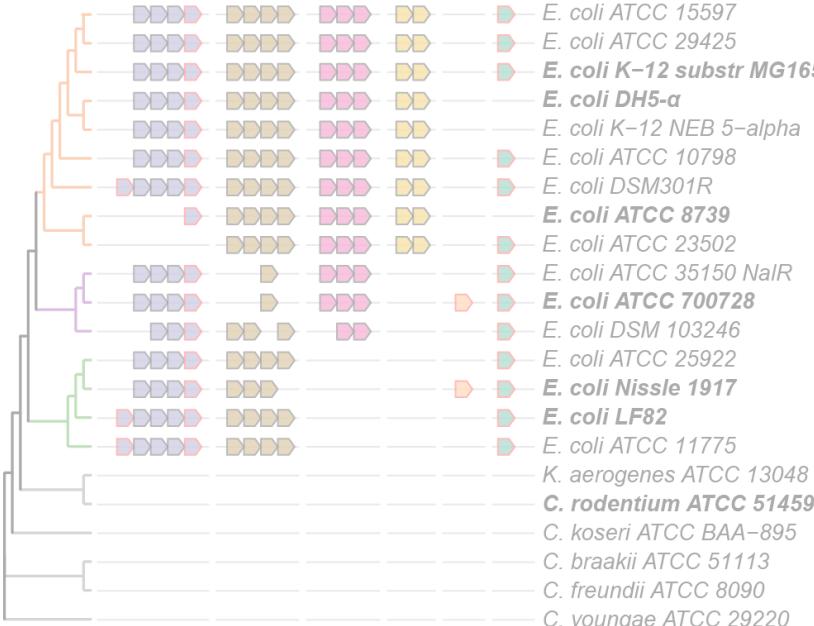


HCT-15 colonic epithelial cells
(MyD88 deficient, no flagellin response)



Bioactivity prioritizations in IBD are supported by external validation

Chaperone-usher pilins
(mostly Type I and P
pilins) were prioritized
during IBD dysbiosis



HCT-15 colonic epithelial cells
(MyD88 deficient, no flagellin response)



Summary

- Introduction to functional profiling
- Tiered search with HUMAnN
- Core functions of the human microbiome
- Contributional diversity
- What's new in bioBakery 3.0?
- Prioritizing bioactive gene families with MetaWIBELE



Microbial community functional profiling

(for real this time)



Summary

- Metatranscriptomics (MTX)
 - Functional profiling of MTX data
 - Relationships between functional potential and activity
 - Detecting under- and over-expression
 - Statistical models for MTX analysis
- Metabolomics (MBX)
 - Sequencing is easy... everything else is hard
 - Metabolomics approaches
 - Identifying compounds
 - Bioactivity prediction with MACARRoN

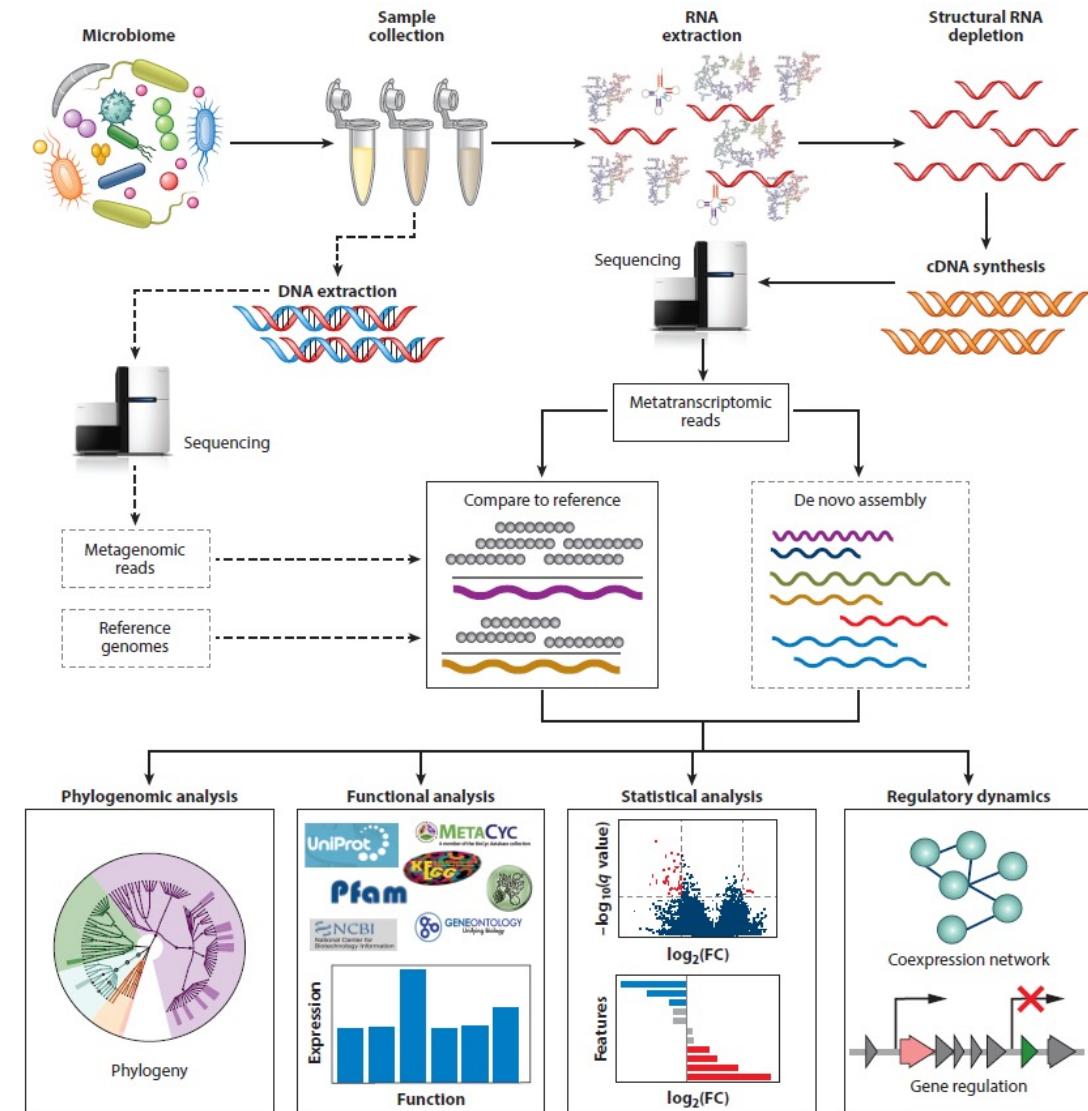


Metatranscriptomics (MTX) overview



Annual Review of Biomedical Data Science Metatranscriptomics for the Human Microbiome and Microbial Community Functional Profiling

Yancong Zhang,^{1,2,*} Kelsey N. Thompson,^{1,2,*}
Tobyn Branck,^{1,2,3} Yan Yan,^{1,2} Long H. Nguyen,^{1,4,5}
Eric A. Franzosa,^{1,2,†} and Curtis Huttenhower^{1,2,6,†}

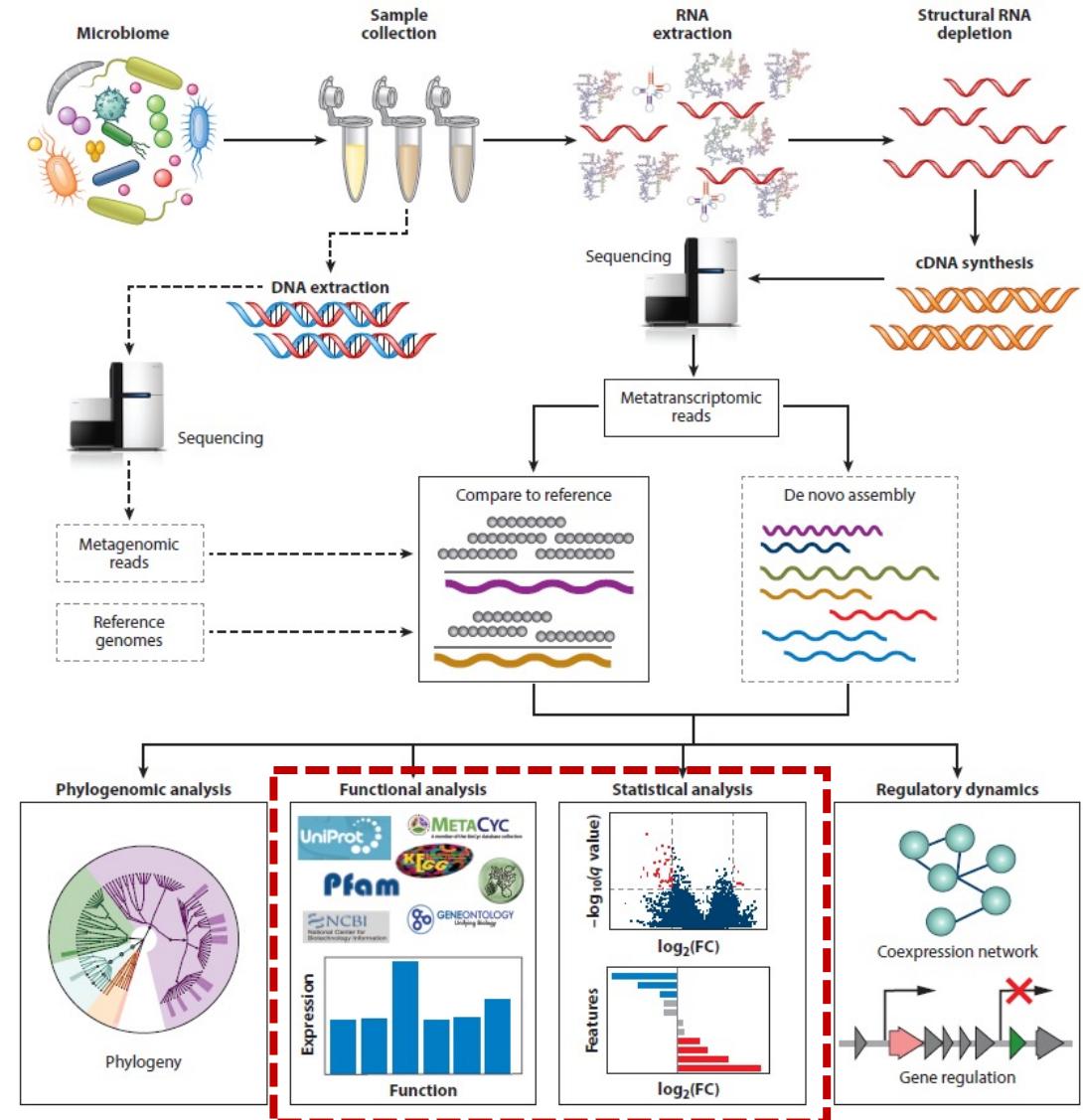
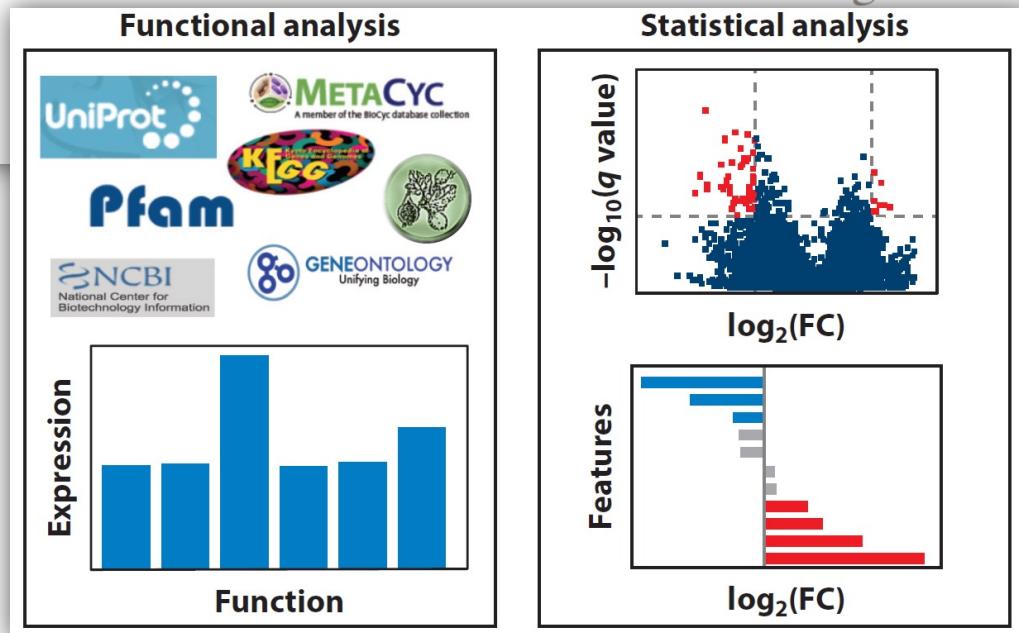




Metatranscriptomics (MTX) overview

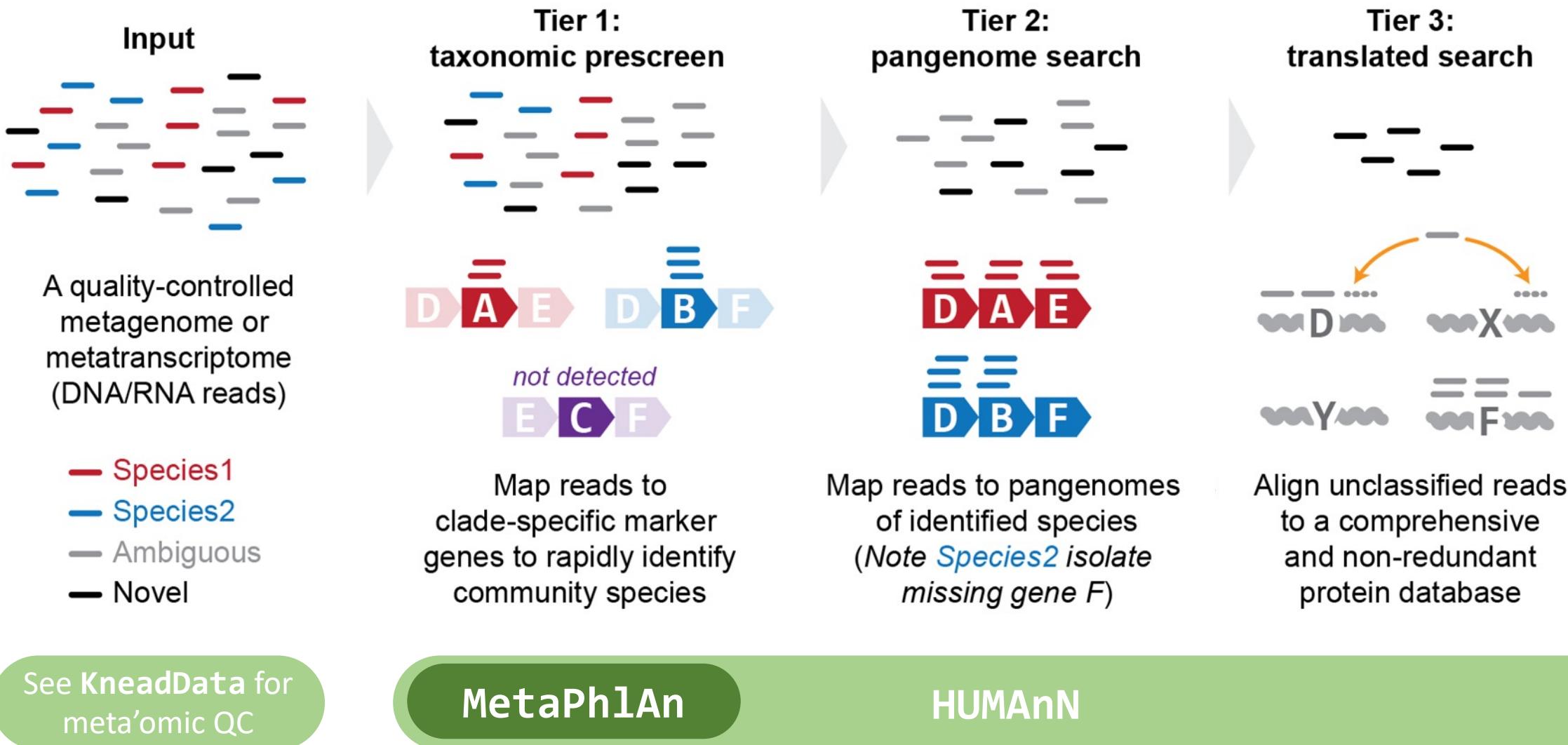


Annual Review of Biomedical Data Science Metatranscriptomics for the Human Microbiome and Microbial Community Functional Profiling



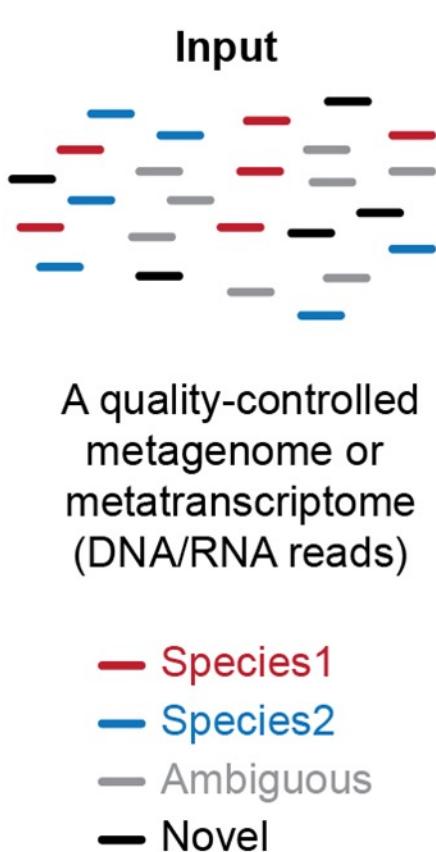


HUMAnN works on metatranscriptomes!





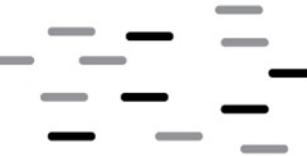
HUMAnN works on metatranscriptomes!



Tier 1:
taxonomic prescreen

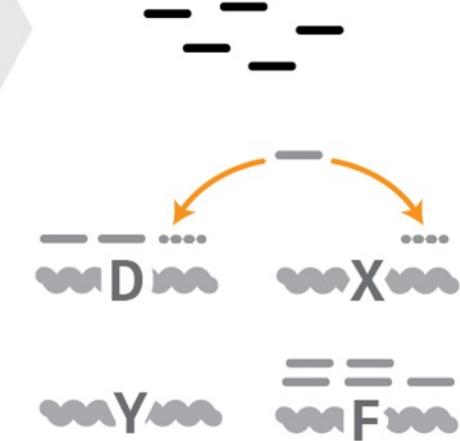
Override
with
taxonomic
profile
from MGX
(if available)

Tier 2:
pangenome search



Map reads to pangenomes
of identified species
(Note *Species2* isolate
missing gene F)

Tier 3:
translated search



Align unclassified reads
to a comprehensive
and non-redundant
protein database

See KneadData for
meta'omic QC

MetaPhlAn

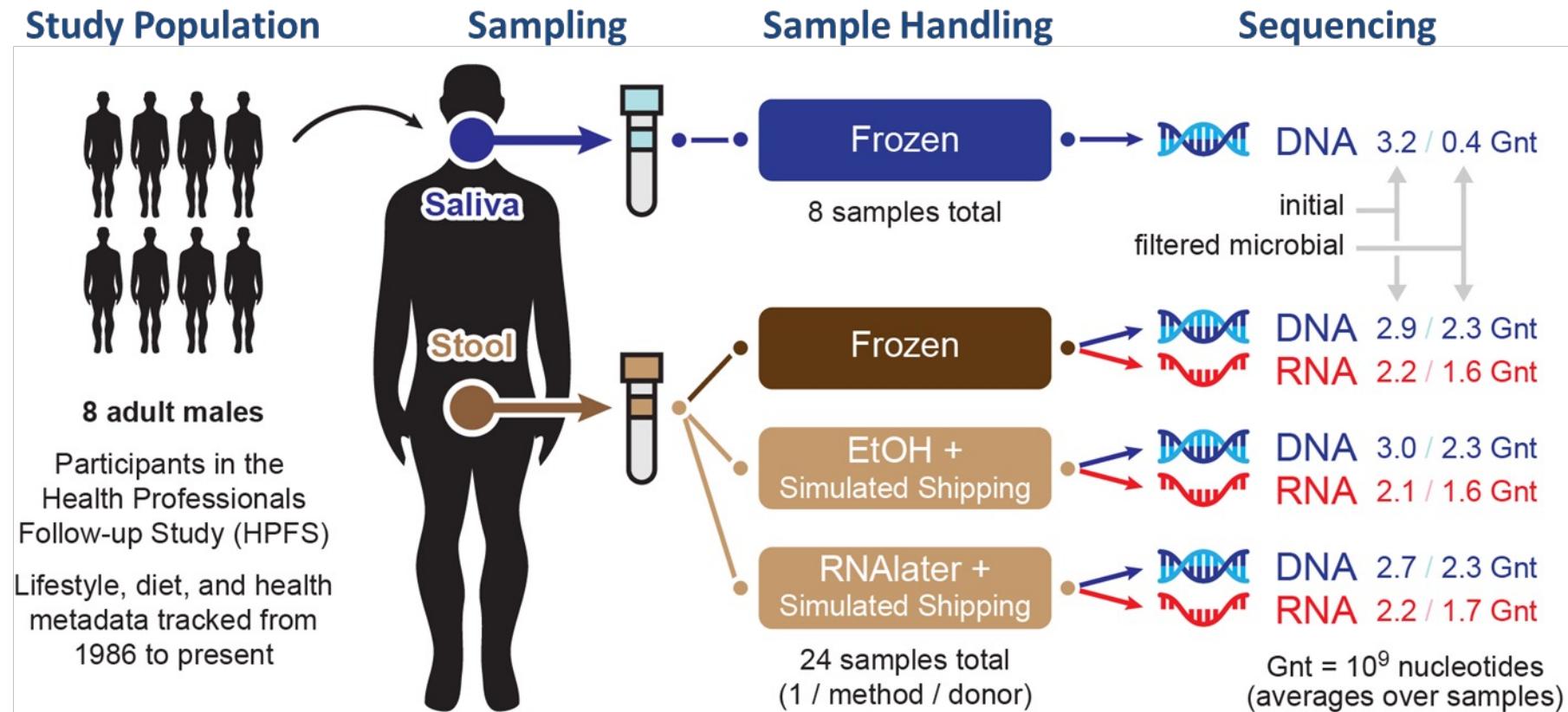
HUMAnN



Piloting metatranscriptomics



The HPFS pilot: self-collected MGX/MTX samples

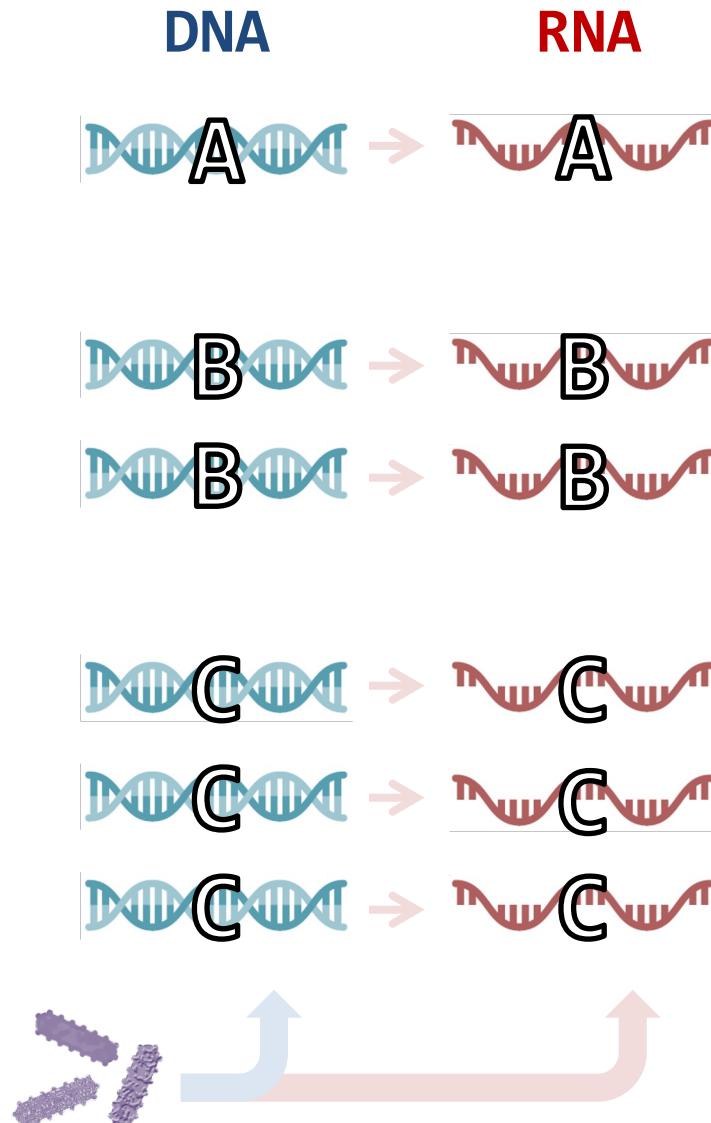


8 self-collected stool samples were mock-shipped with 3 different preservation methods
Then they were subjected to paired MGX and MTX sequencing

Franzosa et al, *Proc Natl Acad Sci U S A*, 2014



Functional potential (MGX/DNA) vs. activity (MTX/RNA)



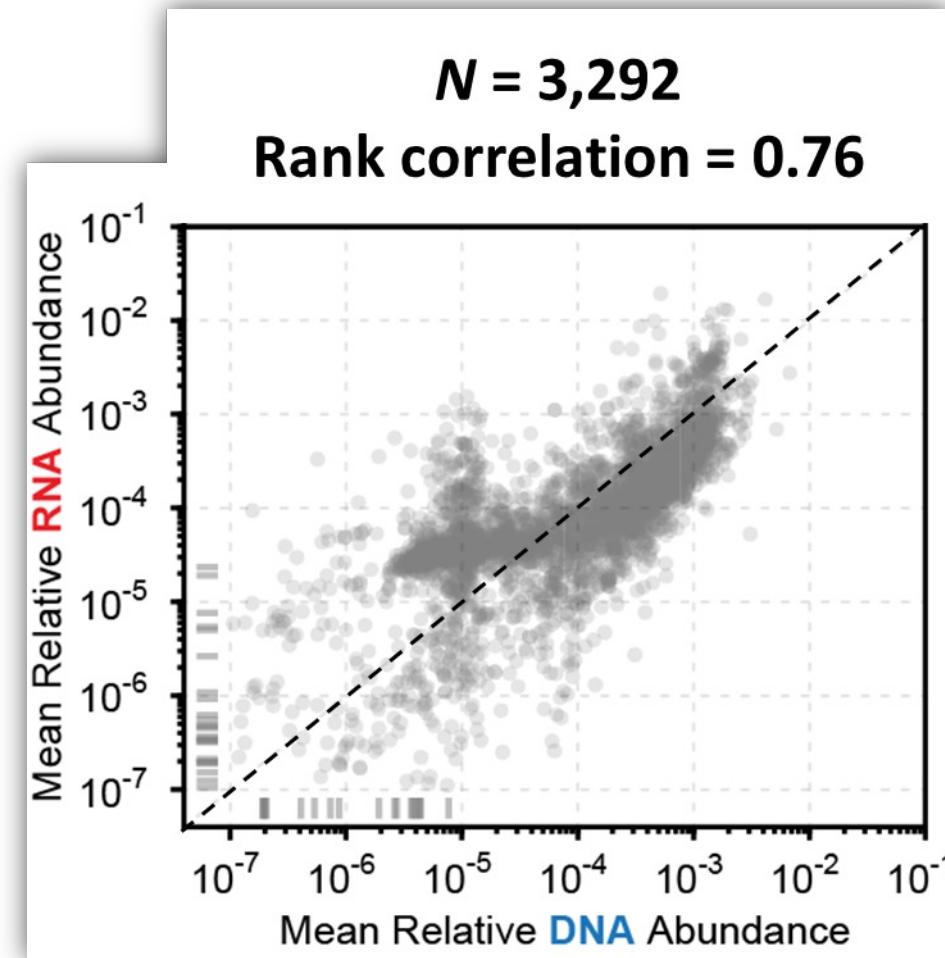
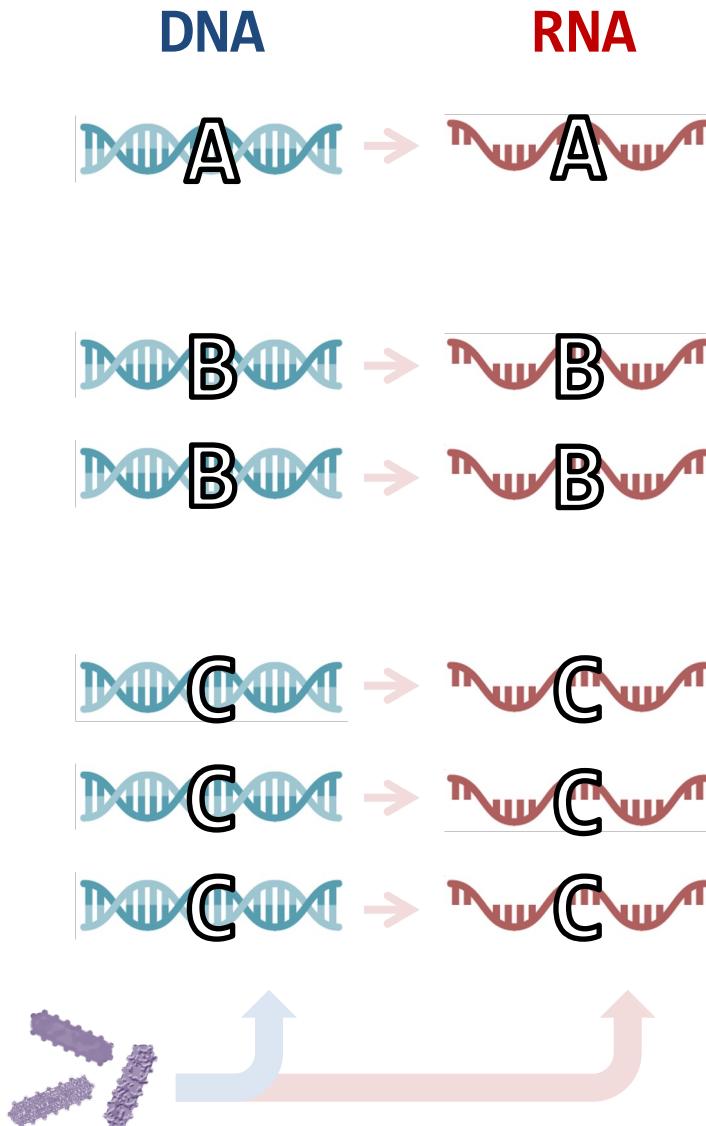
DNA data (metagenomics)
describes functional *potential*

RNA data (metatranscriptomics)
describes functional *activity*

Does *potential* = *activity*
in the human gut?



Potential and activity are strongly correlated!





Potential and activity are strongly correlated!

DNA

RNA



$N = 3,292$

~~Rank correlation = 0.76~~

all else being equal...

more/less gene copies \Rightarrow more/less transcripts



beware that...

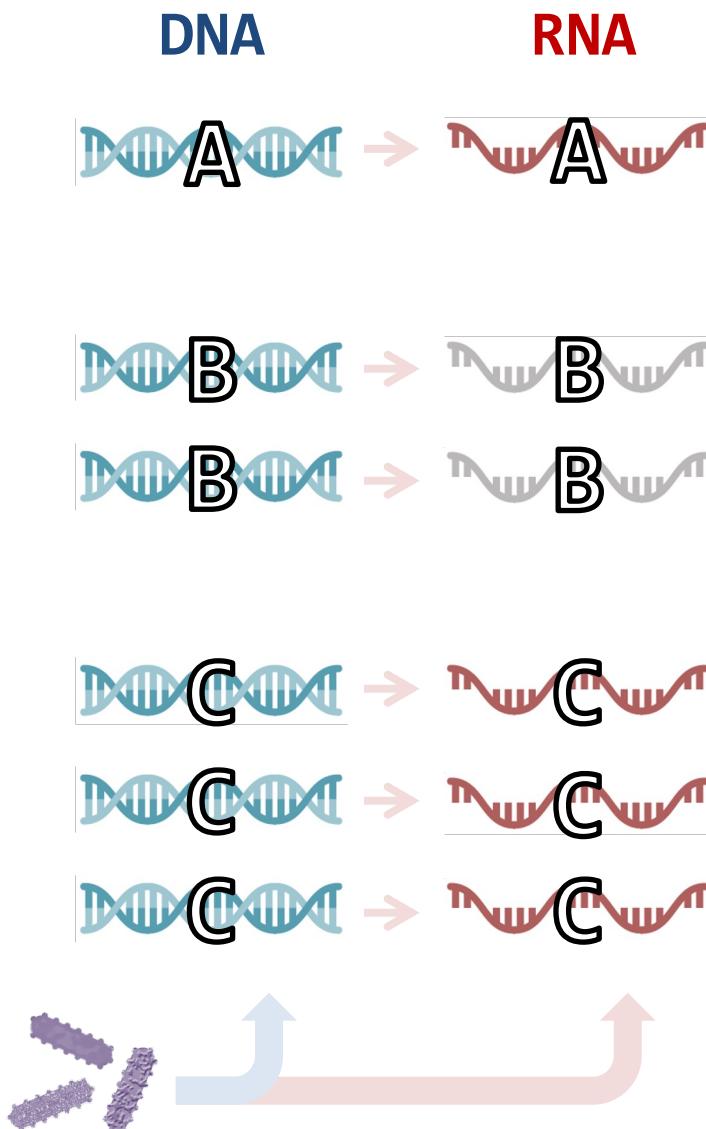
more/less transcripts $\not\Rightarrow$ differential expression
(could be changes in composition, i.e. gene copies)



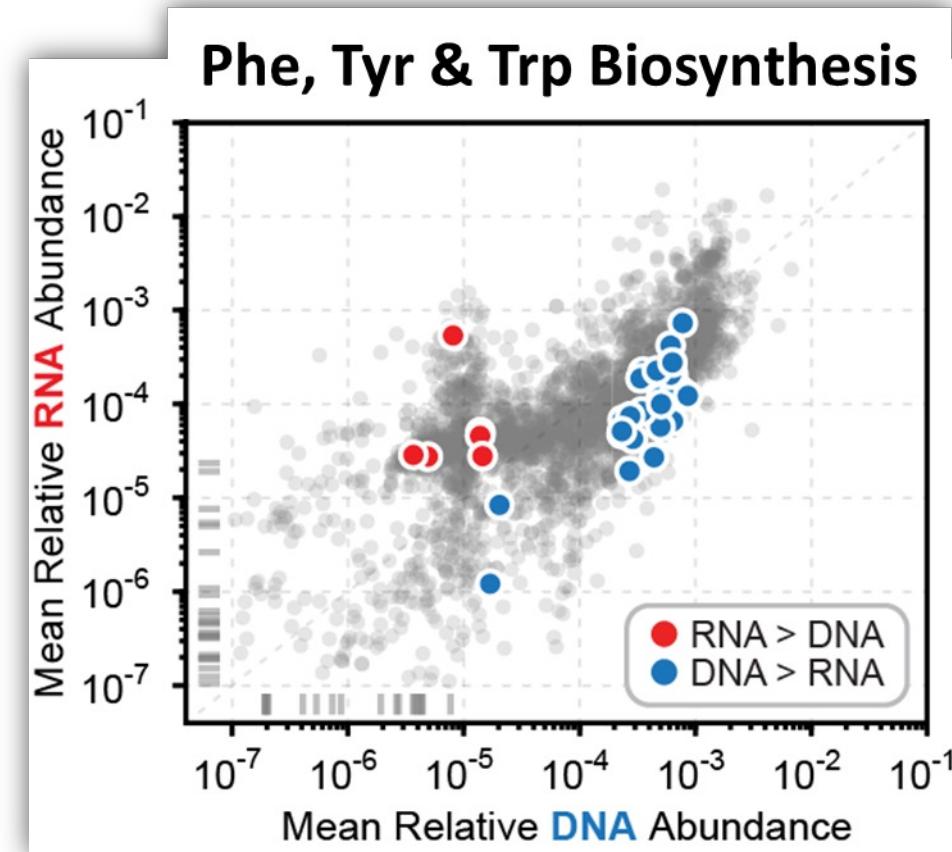
$10^{-7} \quad 10^{-6} \quad 10^{-5} \quad 10^{-4} \quad 10^{-3} \quad 10^{-2} \quad 10^{-1}$
Mean Relative DNA Abundance



Some functions are “under-abundant” in MTX

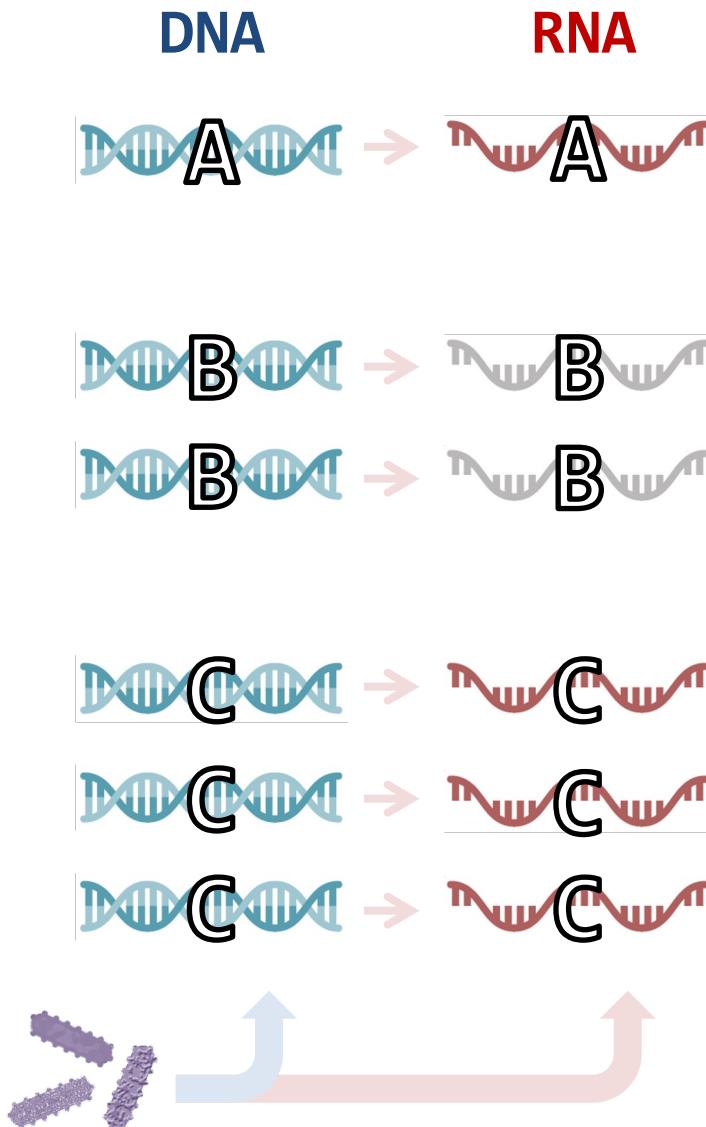


Low RNA abundance relative to DNA abundance

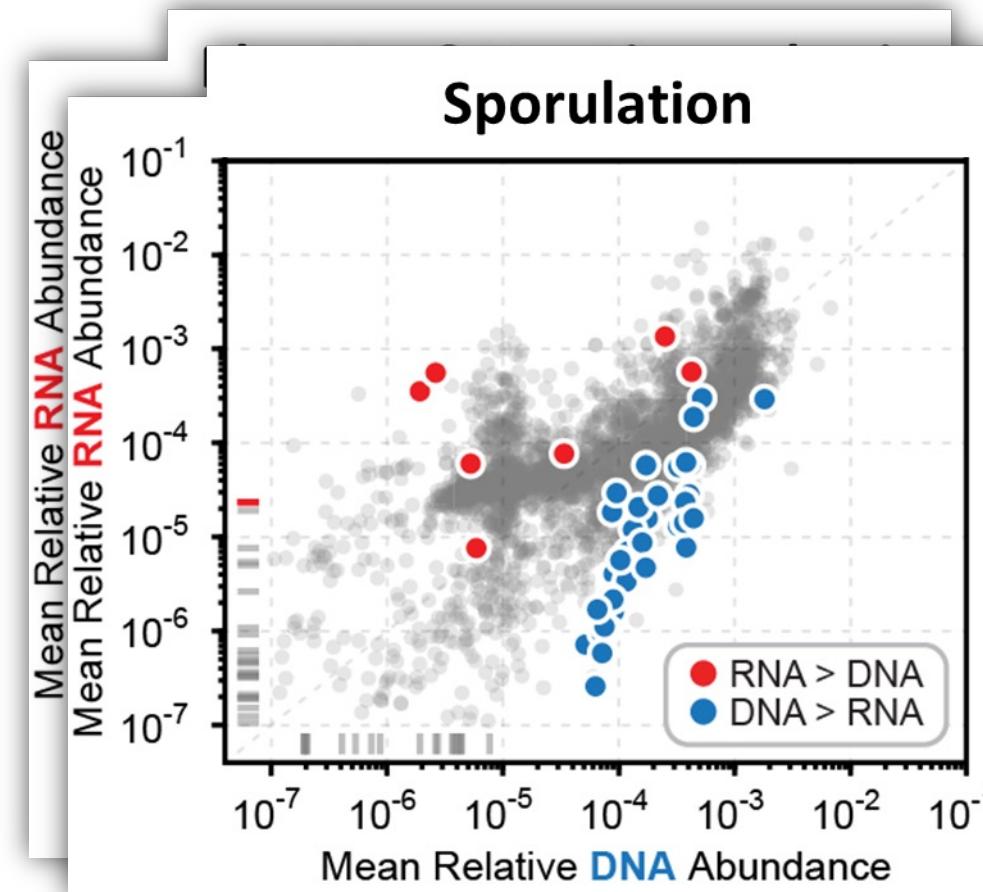




Some functions are “under-abundant” in MTX

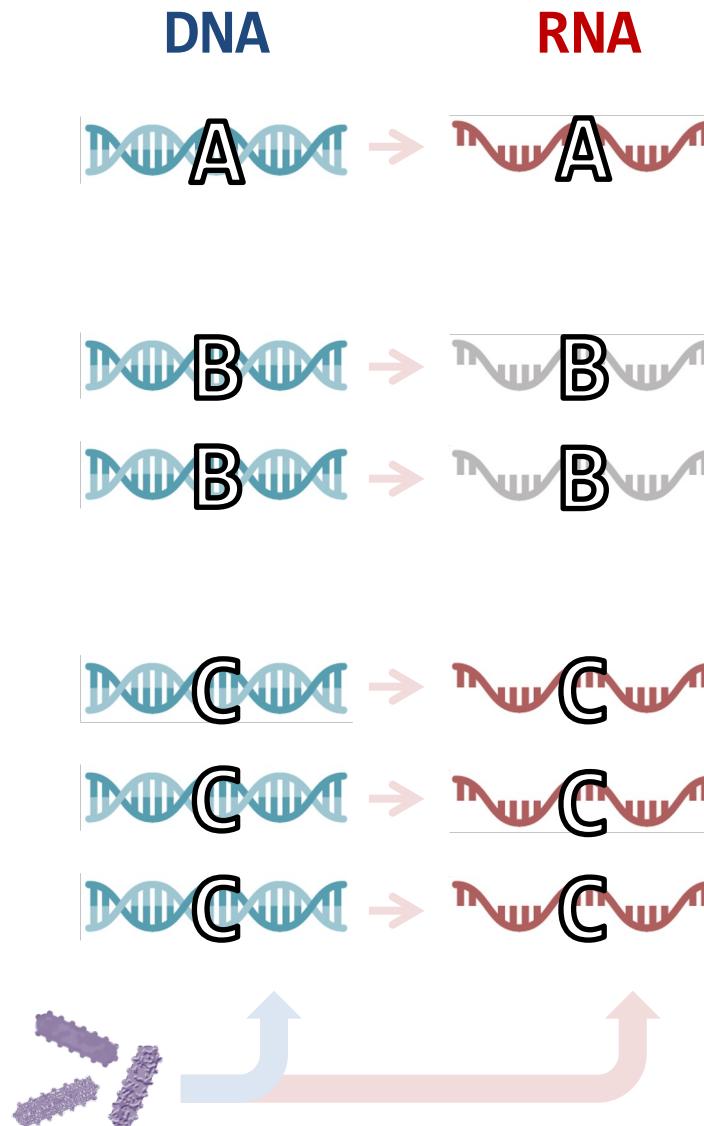


Low RNA abundance relative to DNA abundance

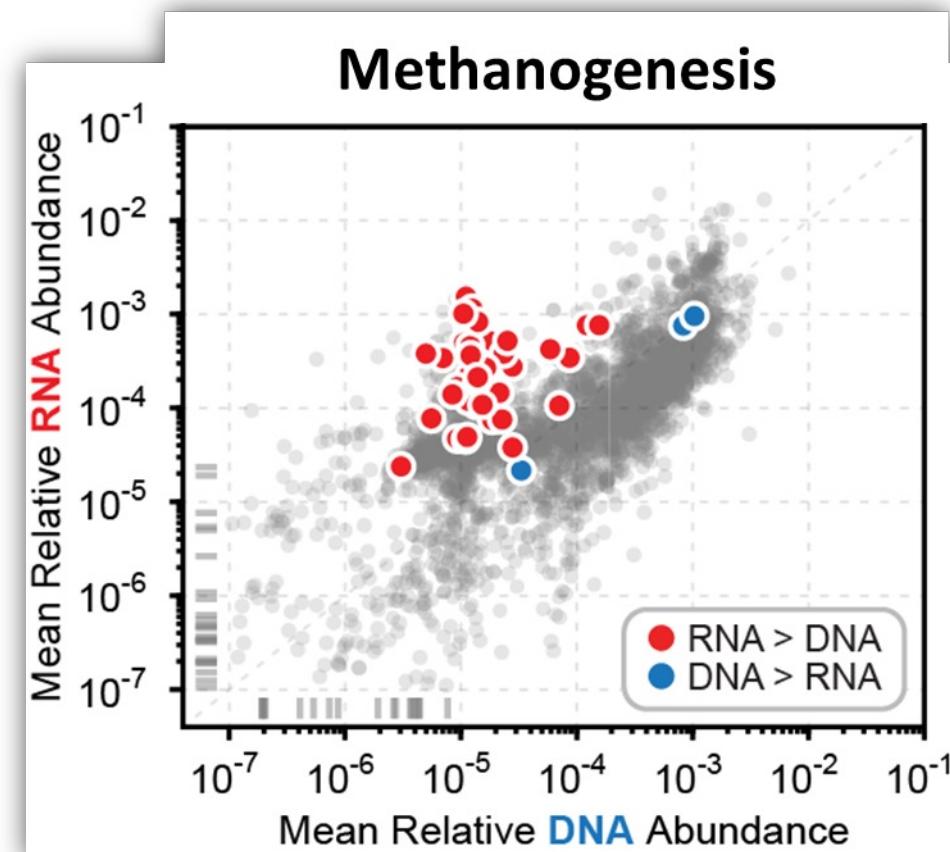




Other functions are “over-abundant” in MTX

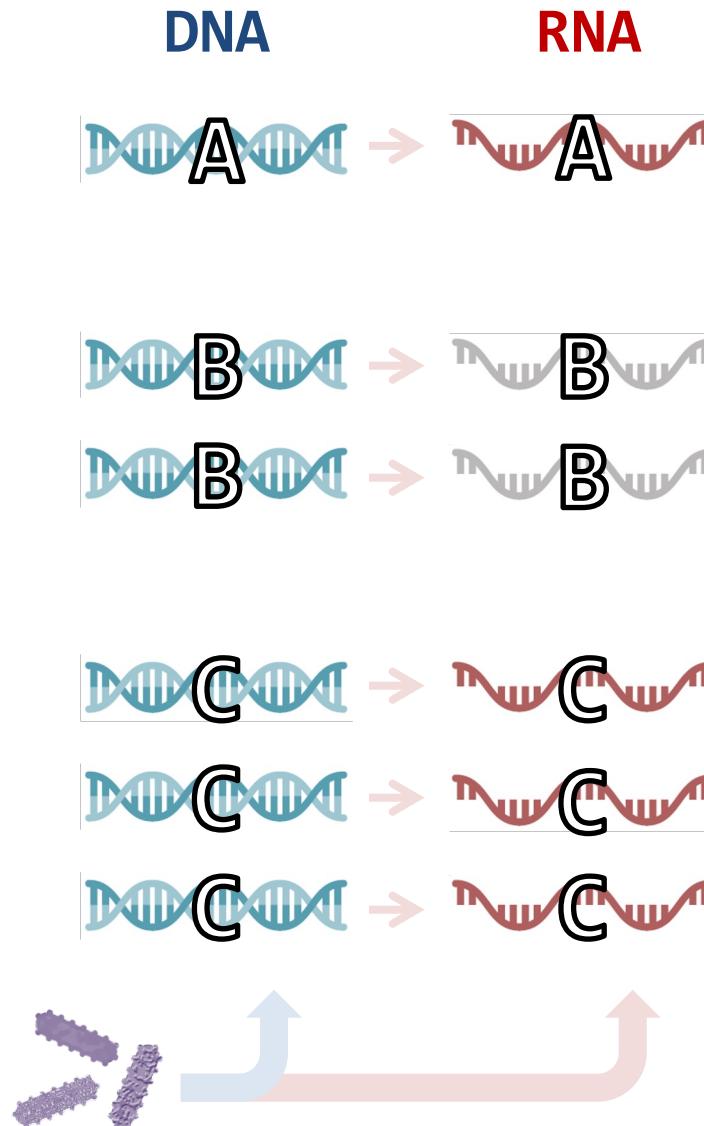


High RNA abundance relative to DNA abundance

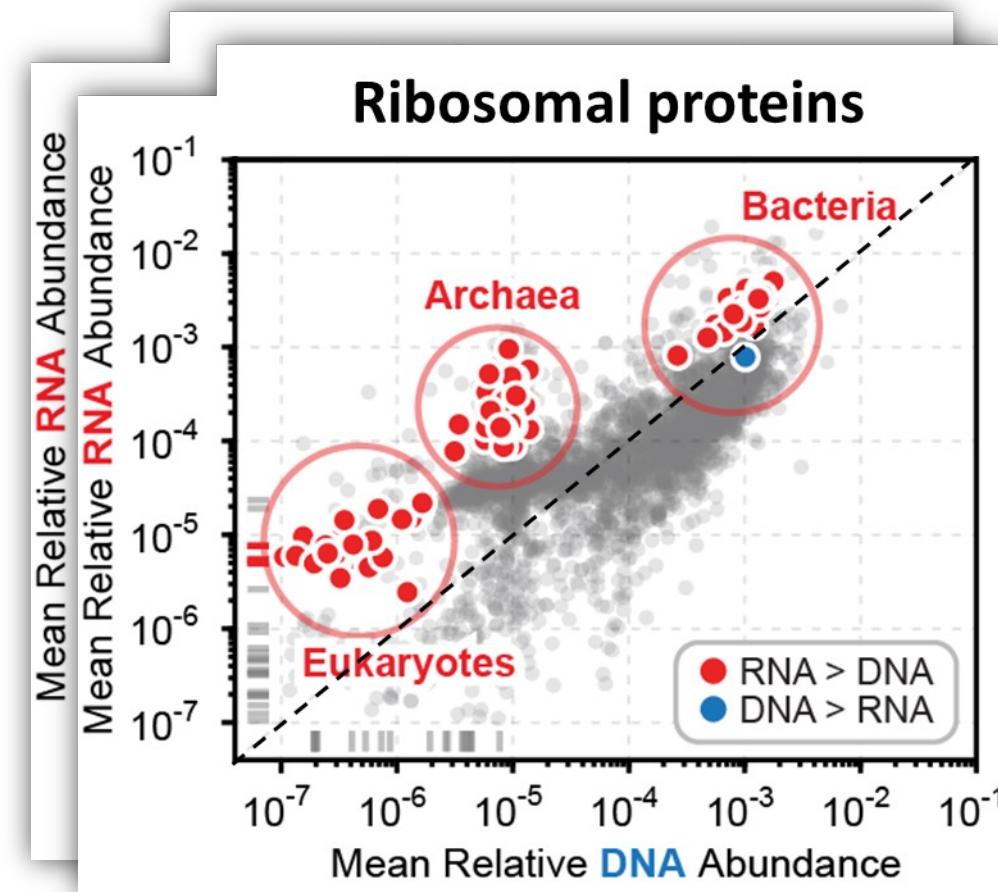




Other functions are “over-abundant” in MTX

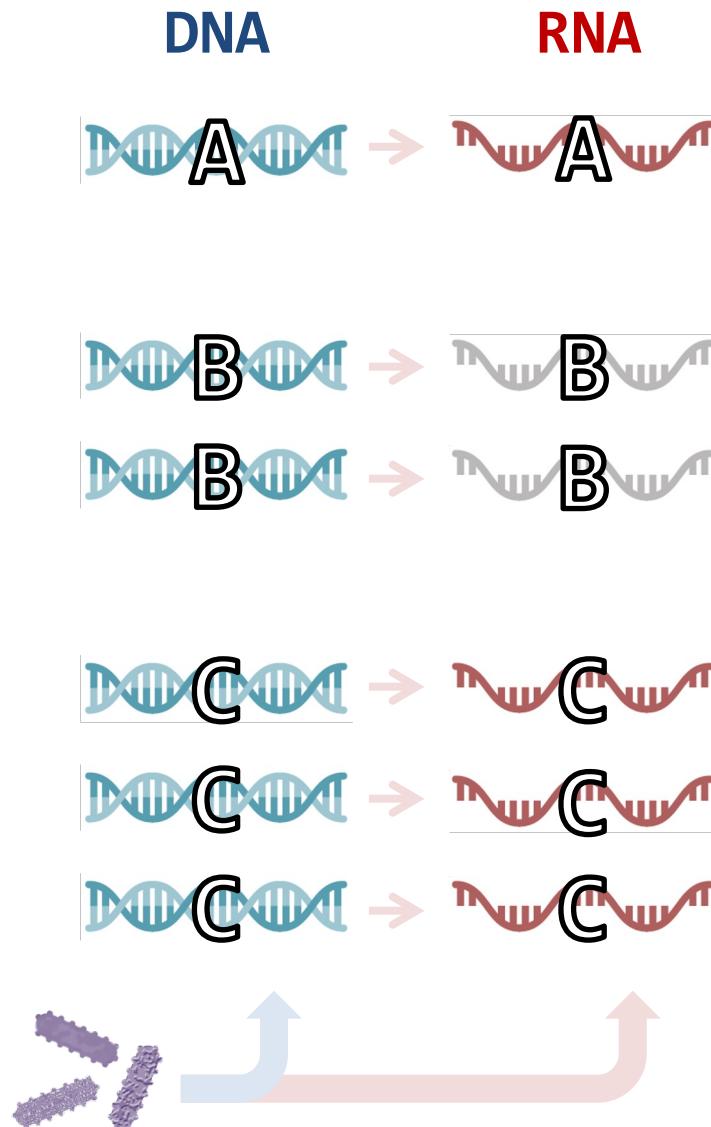


High RNA abundance relative to DNA abundance

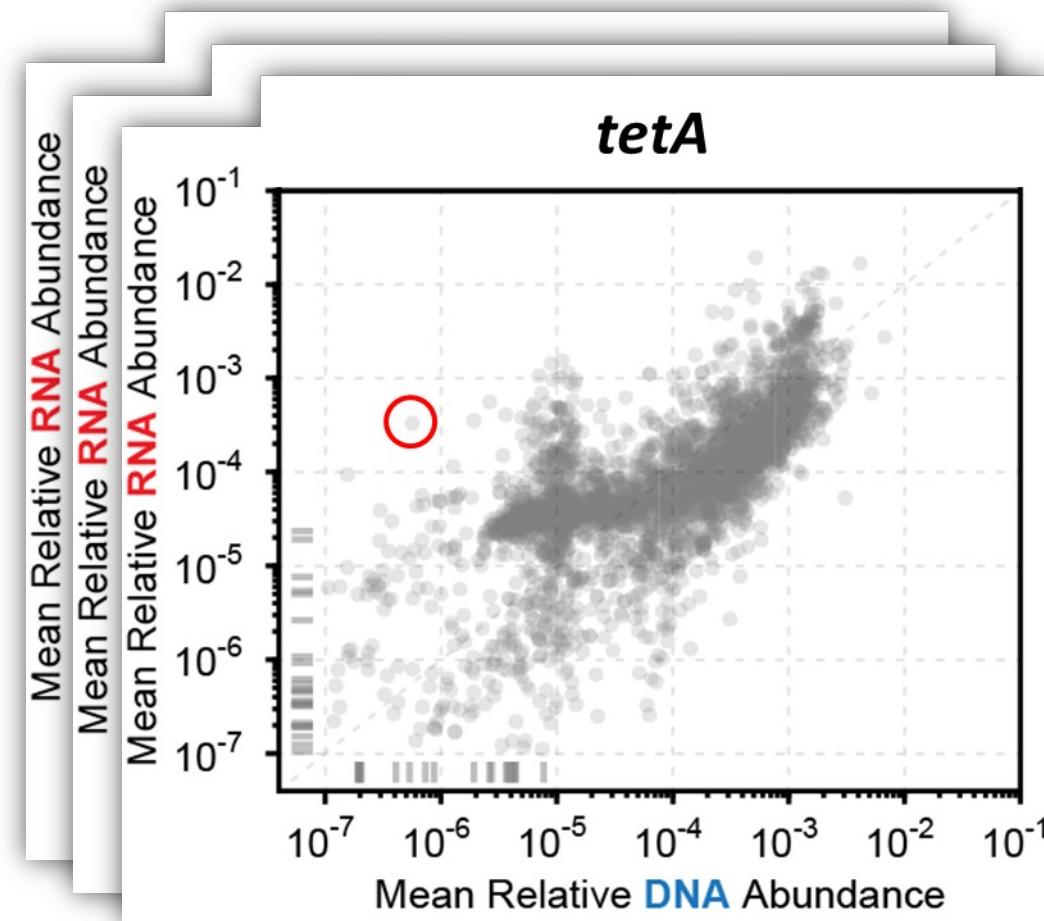




Other functions are “over-abundant” in MTX

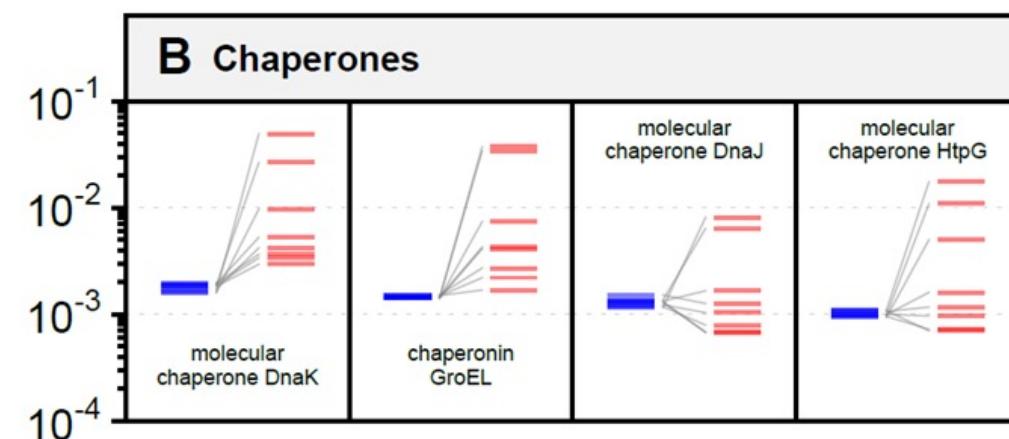
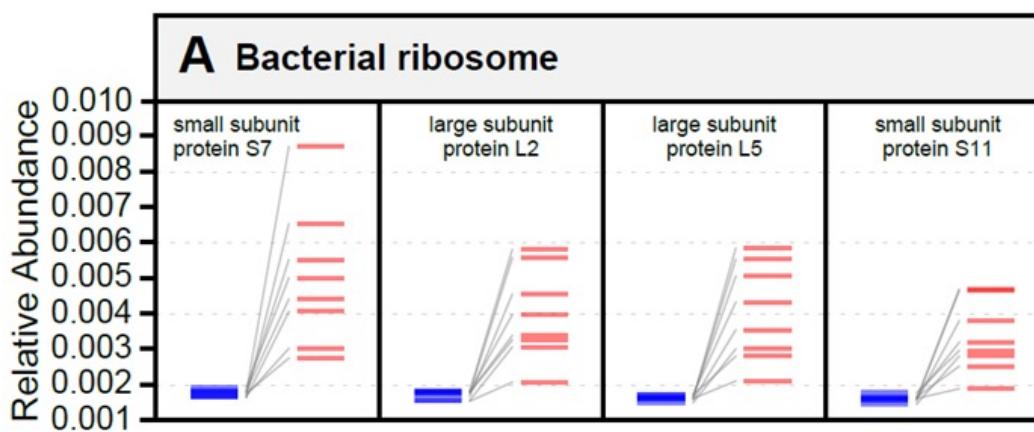
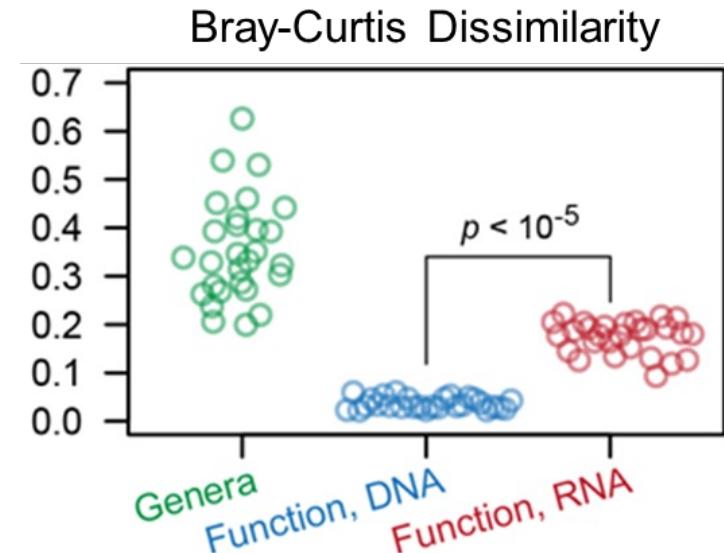


High RNA abundance relative to DNA abundance



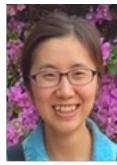


Activity is more personalized than potential





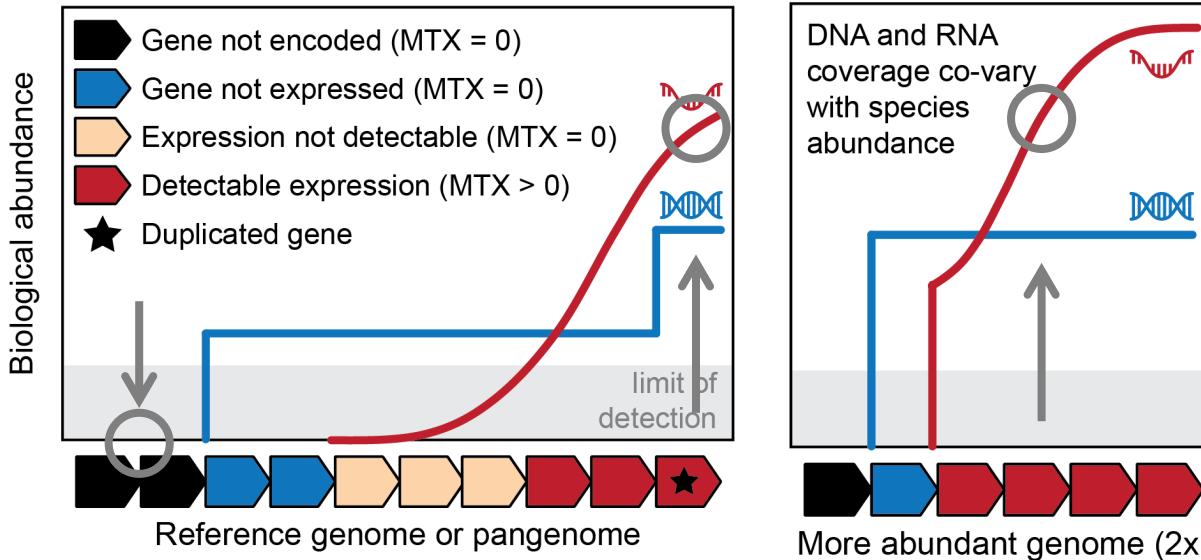
Statistics for metatranscriptomics (MTX)



Yancong
Zhang



Stats challenges in the analysis of MTX data



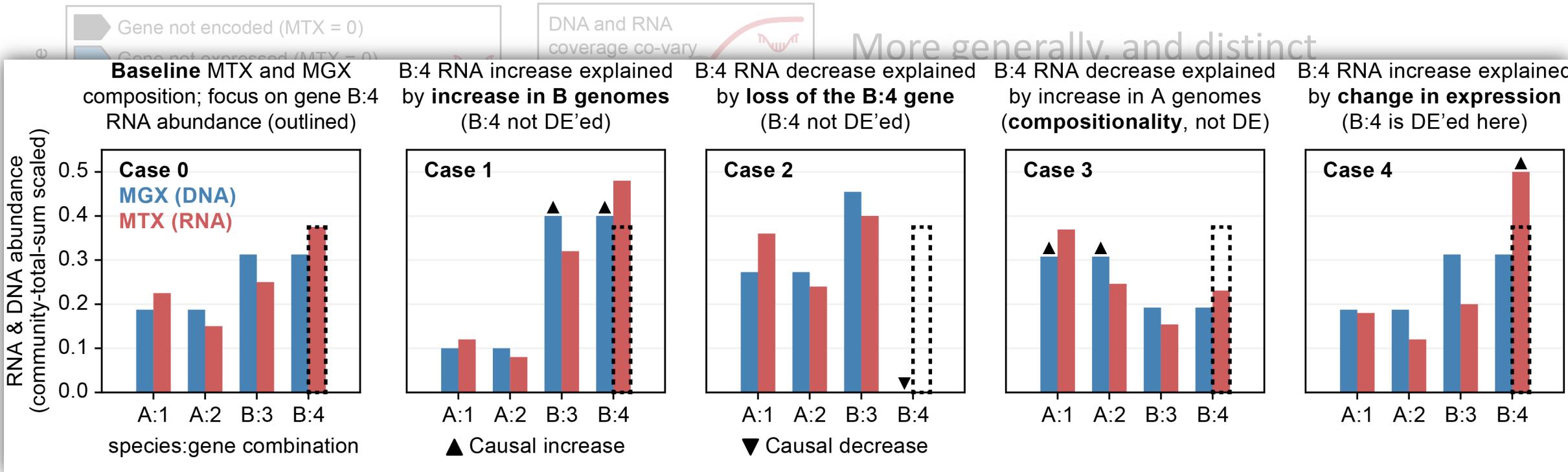
More generally, and distinct from single-species RNA-seq, MTX copy number depends on MGX (gene) copy number

Technical zeros are common due to dynamic range : seq depth interaction

Biological zeros are common from missing species or gene deletion



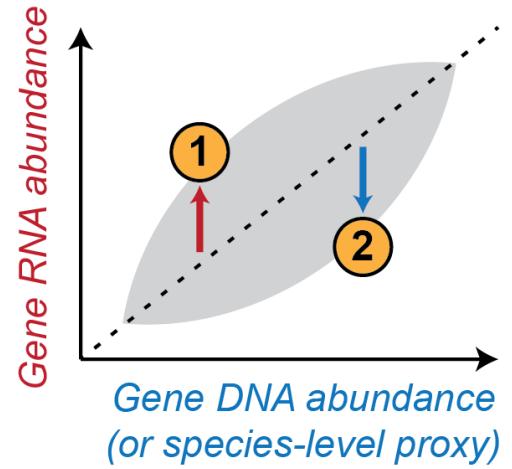
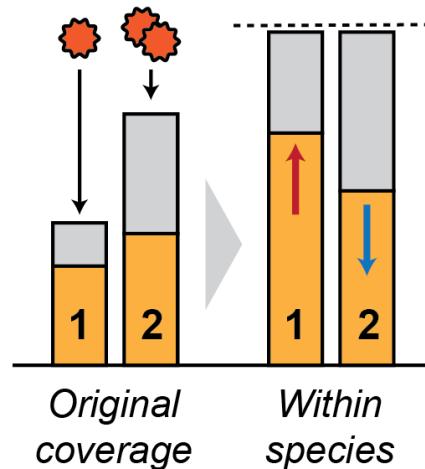
Stats challenges in the analysis of MTX data



As a consequence, many phenomena that lead to differential abundance in MTX data are not related to differential regulation of microbial genes!



Accounting for “relative expression” (concepts)



Family 1: normalize
RNA within each
community species
 \Rightarrow
(essentially)
 N separate single-
species RNA-seq
datasets

Family 2: adjust
RNA abundance for
gene copy number
(or a proxy, e.g.
species abundance)
using a ratio or
residual



Six formal models of relative expression...

Model number: name	Model formulation	Model description
M1: Naïve RNA	$C(f_{RNA}) \sim p$	The community total-sum-scaled abundance $C(x)$ of a feature f 's RNA count (f_{RNA}) is modeled as a function of a sample phenotype/property (p).
✗ M2: Within-taxon RNA	$T(f_{RNA}) \sim p$	f_{RNA} is subjected to within-taxon sum-scaling, $T(x)$, i.e. normalizing against the total RNA pool of f 's source taxon.
✗ M3: Taxon-RNA covariate	$T(f_{RNA}) \sim Tax_{RNA}(f) + p$	f_{RNA} is scaled within-taxon as in M2; an RNA-level estimate of f 's source taxon abundance, $Tax_{RNA}(f)$, is included as a covariate.
+ M4: RNA/DNA ratio	$C(f_{RNA})/C(f_{DNA}) \sim p$	The ratio of f 's total-sum scaled RNA and DNA abundances (the “relative expression ratio”) is modeled as a function of p .
* M5: Taxon-DNA covariate	$C(f_{RNA}) \sim Tax_{DNA}(f) + p$	f_{RNA} is total-sum scaled as in M1; a DNA-level estimate of f 's source taxon abundance, $Tax_{DNA}(f)$, is included as a covariate.
+ M6: Feature-DNA covariate	$C(f_{RNA}) \sim C(f_{DNA}) + p$	f_{RNA} is total-sum scaled as in M1; a total-sum scaled estimate of f 's DNA-level abundance is included as a covariate.

✗ = The model requires a mapping of features to taxa + = the model requires paired MGX (DNA) data * = both are required



...and three modes for pre-filtering technical zeros

Community RNA + DNA
for feature f over 6 samples

f_{RNA}	1	0	1	0	2	0
f_{DNA}	0	0	1	0	1	1

Lenient filtering: Treat all zeroes as informative but ignore features that are never seen at the RNA or DNA level.

exclude 2 samples

f_{RNA}	1	0	1	0	2	0
f_{DNA}	0	0	1	0	1	1

Semi-strict filtering: Treat zero values as informative if $\max(f_{\text{RNA}}, f_{\text{DNA}}) > 0$; exclude samples with $\max(f_{\text{RNA}}, f_{\text{DNA}}) = 0$.

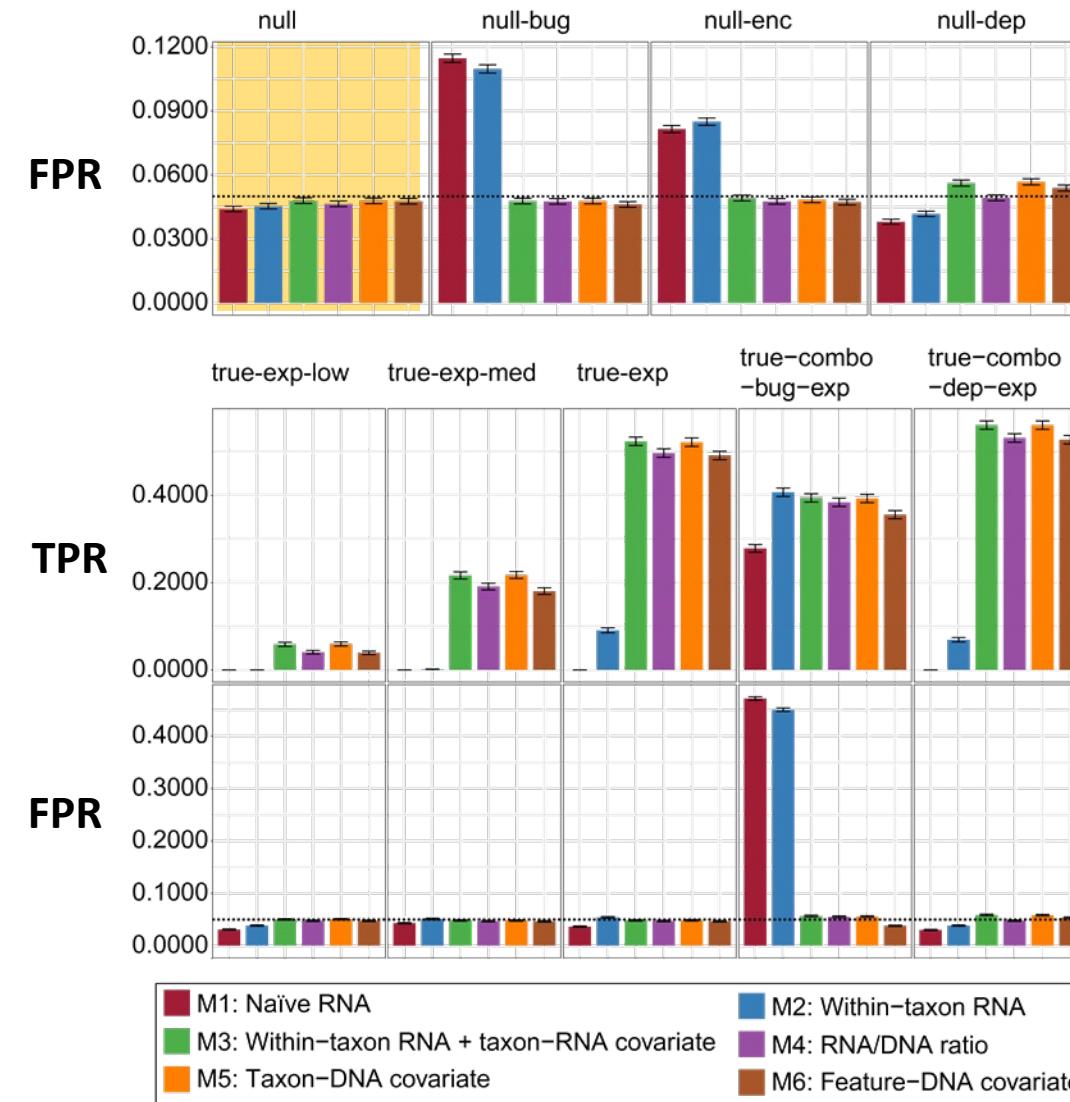
exclude 4 samples

f_{RNA}	1	0	1	0	2	0
f_{DNA}	0	0	1	0	1	1

Strict filtering: Treat all zeroes as likely technical and uninformative; exclude samples with $\min(f_{\text{RNA}}, f_{\text{DNA}}) = 0$.



Model performance on synthetic MTX datasets



- **(Based on strict filtering of zeros!)**
- All methods report OK FPR with no signal
- M1 (naïve) and M2 (within-species RNA) confuse bug abundance changes for DE
- M3 (within-species RNA with total RNA covariate) has better FPR and good TPR (no DNA needed)
- DNA-aware methods are roughly equivalent here
- Note: Strict filtering of MTX zeros prevents models M1-3 & 5 from confusing gene loss for DE



Model performance on synthetic MTX datasets



- **(Based on strict filtering of zeros!)**
- All methods report OK FPR with no signal
- M1 (naïve) and M2 (within-species RNA) confuse bug abundance changes for DE
- M3 (within-species RNA with total RNA covariate) has better FPR and good TPR (no DNA needed)
- DNA-aware methods are roughly equivalent here
- Note: Strict filtering of MTX zeros prevents models M1-3 & 5 from confusing gene loss for DE



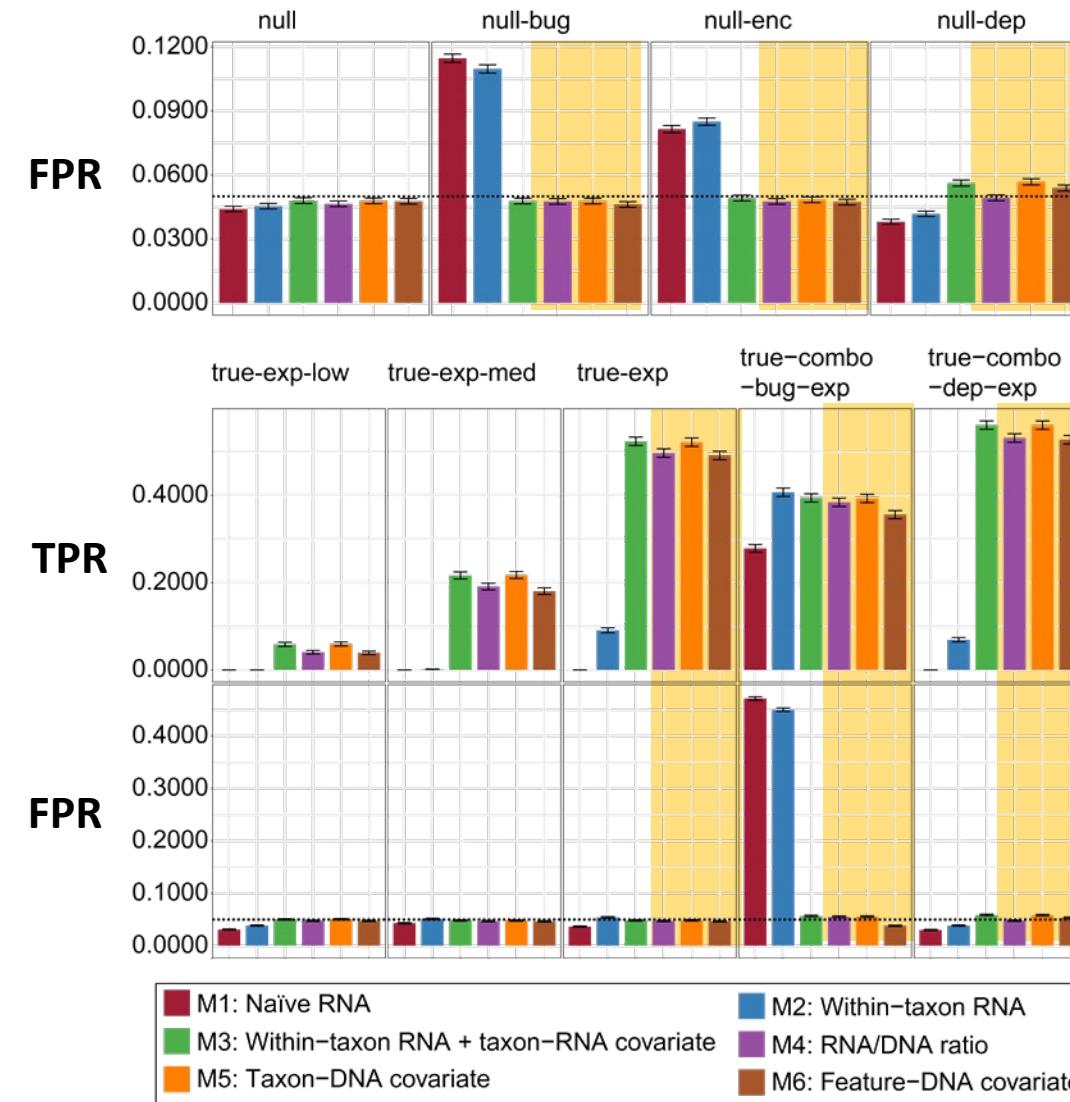
Model performance on synthetic MTX datasets



- (*Based on strict filtering of zeros!*)
- All methods report OK FPR with no signal
- M1 (naïve) and M2 (within-species RNA) confuse bug abundance changes for DE
- M3 (within-species RNA with total RNA covariate) has better FPR and good TPR (no DNA needed)
- DNA-aware methods are roughly equivalent here
- Note: Strict filtering of MTX zeros prevents models M1-3 & 5 from confusing gene loss for DE



Model performance on synthetic MTX datasets



- (*Based on strict filtering of zeros!*)
- All methods report OK FPR with no signal
- M1 (naïve) and M2 (within-species RNA) confuse bug abundance changes for DE
- M3 (within-species RNA with total RNA covariate) has better FPR and good TPR (no DNA needed)
- DNA-aware methods are roughly equivalent here
- Note: Strict filtering of MTX zeros prevents models M1-3 & 5 from confusing gene loss for DE



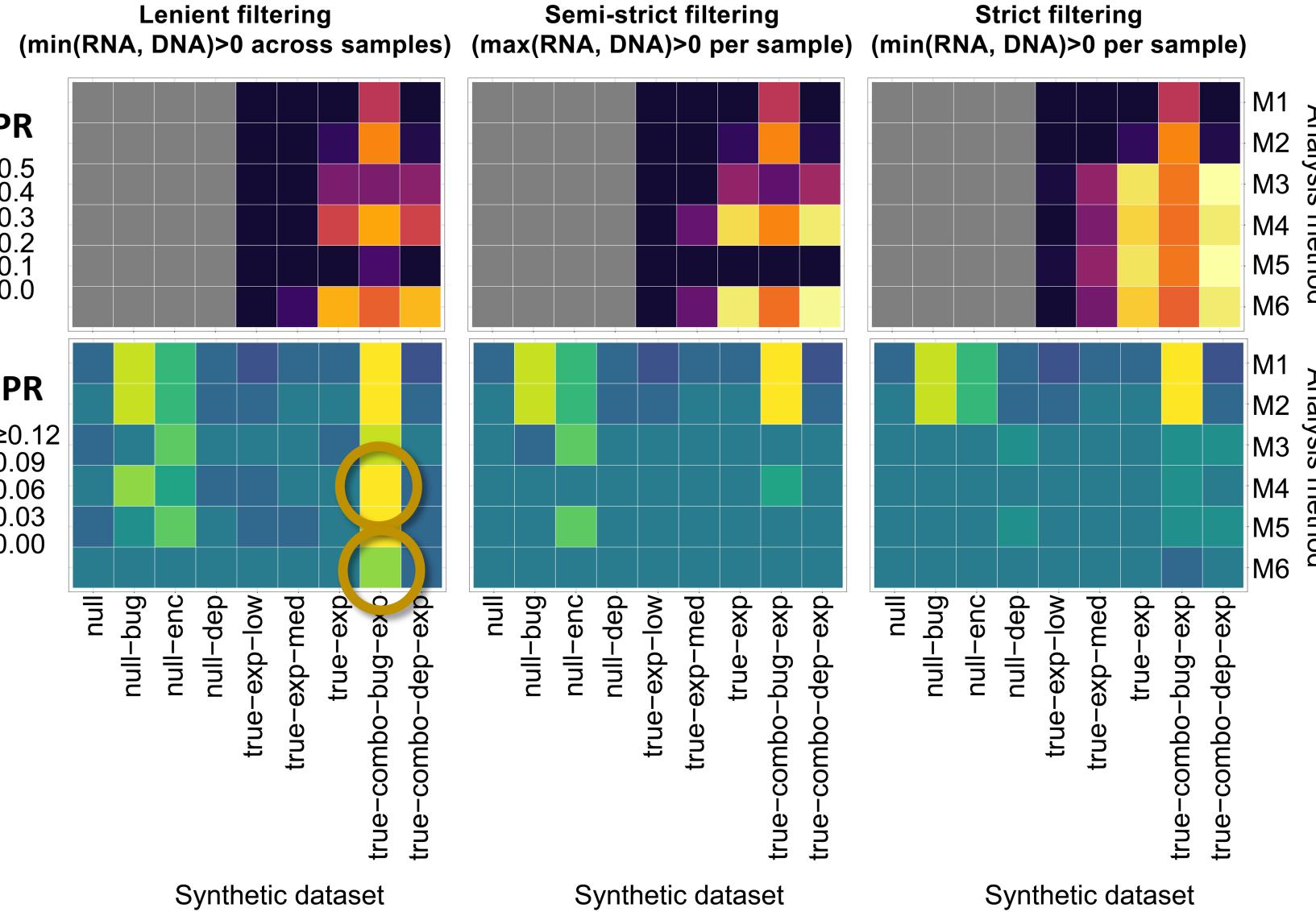
Model performance on synthetic MTX datasets



- (*Based on strict filtering of zeros!*)
- All methods report OK FPR with no signal
- M1 (naïve) and M2 (within-species RNA) confuse bug abundance changes for DE
- M3 (within-species RNA with total RNA covariate) has better FPR and good TPR (no DNA needed)
- DNA-aware methods are roughly equivalent here
- Note: Strict filtering of MTX zeros prevents models M1-3 & 5 from confusing gene loss for DE



Model performance vs. zero-filtering method



- These are the results from the previous slide.
- M3 (best no-DNA model) struggles in the presence of technical zeros AND calls gene loss as DE.
- M5 (normalizing by species DNA abundance) behaves similarly.
- M6 (gene DNA covariate) is more robust to zeros than M4 (RNA/DNA ratio) given strong confounding.



Current recommendations for MTX stats

Model number: name	Model formulation
M1: Naïve RNA	$C(f_{RNA}) \sim p$
✗ M2: Within-taxon RNA	$T(f_{RNA}) \sim p$
✗ M3: Taxon-RNA covariate	$T(f_{RNA}) \sim Tax_{RNA}(f) + p$
+ M4: RNA/DNA ratio	$C(f_{RNA})/C(f_{DNA}) \sim p$
* M5: Taxon-DNA covariate	$C(f_{RNA}) \sim Tax_{DNA}(f) + p$
+ M6: Feature-DNA covariate	$C(f_{RNA}) \sim C(f_{DNA}) + p$

If you **don't** have paired MGX data, M3 is a reasonable approach if...

- You can assign your transcripts to species
- You ignore all zero values

If you **do** have paired MGX data, M4 and both M6 behave well, with M6 (the gene-DNA covariate) proving more robust to zero values than the RNA/DNA ratio.

M6 is our current default.

✗ = The model requires a mapping of features to taxa

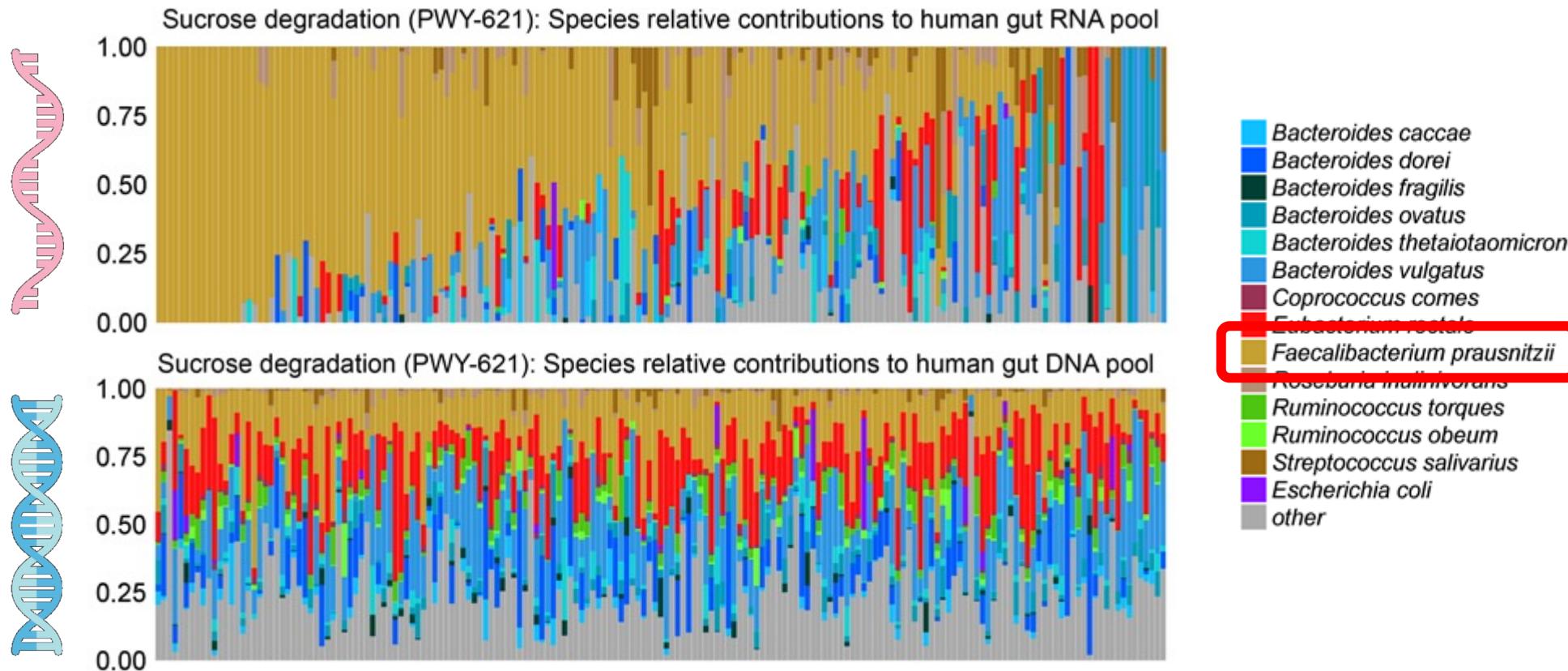
+ = the model requires paired MGX (DNA) data

* = both are required



Aside: Contributional changes in activity

Abu-Ali et al.,
Nat Microbiol, 2017



Sucrose degradation follows a complex attribution pattern across ~200 human gut metagenomes...
...but its expression can be dominated by a single species in paired gut metatranscriptomes!

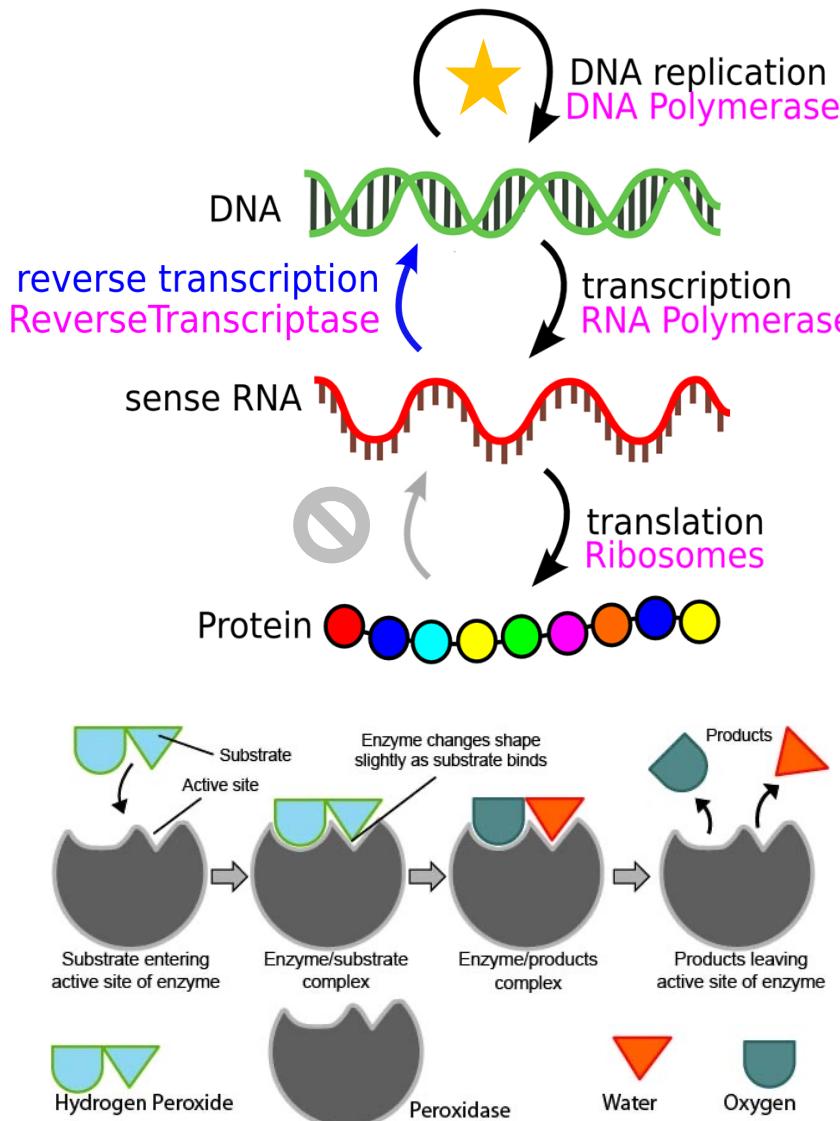
In collaboration with
the STARR Consortium
& HPFS cohort



Metabolomics



Why metabolomics and proteomics are different

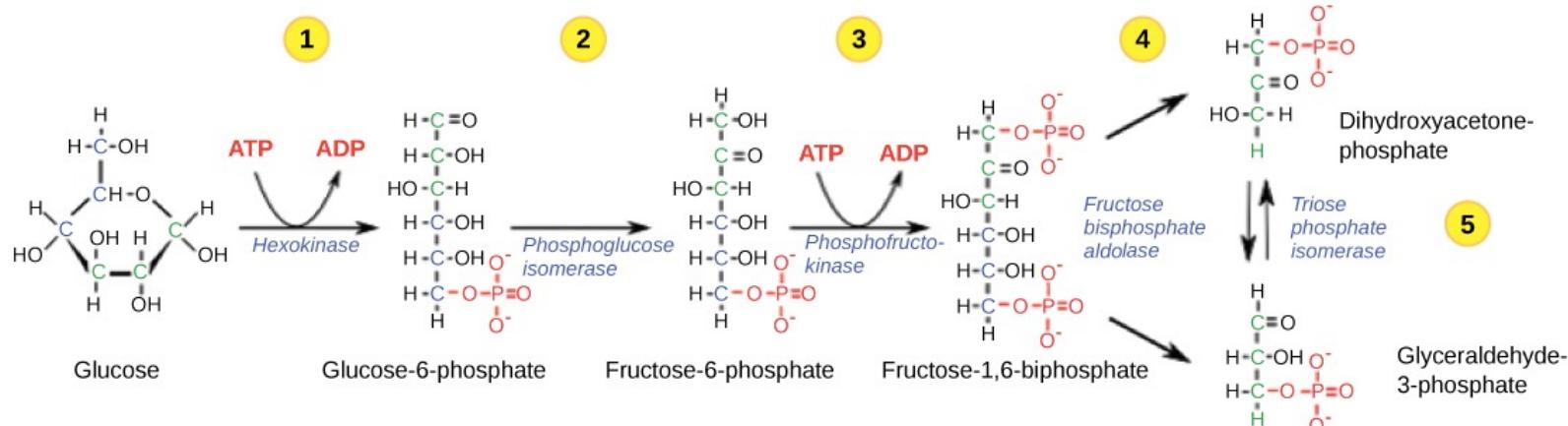
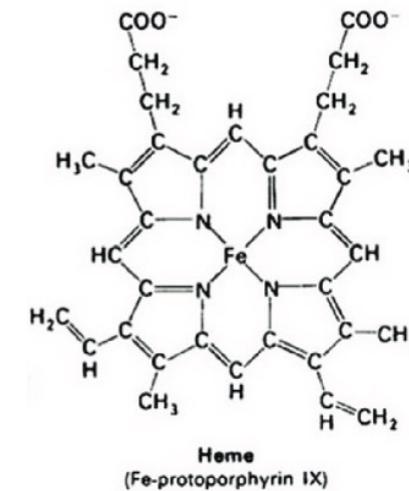


- DNA “wants” to replicate
 - Most sequencing is sequencing by synthesis
- Many problems solved directly by DNA sequencing
 - Genomics, metagenomic profiling
- Other ‘omics problems solved by reducing to DNA seq.
 - Transcriptomics (reverse-transcribe RNA to DNA, then seq.)
 - Chromatin conformation (split and rejoin local DNA, then seq.)
 - Protein interaction (seq. interacting genes in rescued cells)
- No biological “reverse translation” (protein → DNA)
 - Quantifying proteins (proteomics) needs another approach
 - (Aside: we can and will map protein to DNA with e.g. BLAST)
- Metabolites aren’t a part of the central dogma
 - (Though they *are* acted on by proteins)
 - Metabolomics will require different approaches
 - *We’ll start here...*



What is metabolomics?

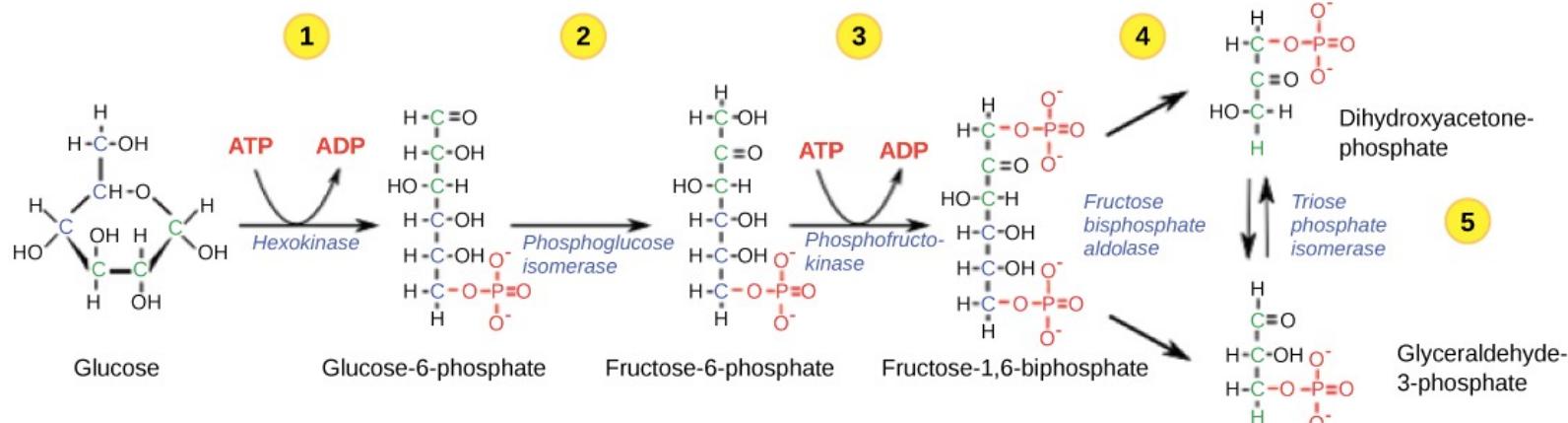
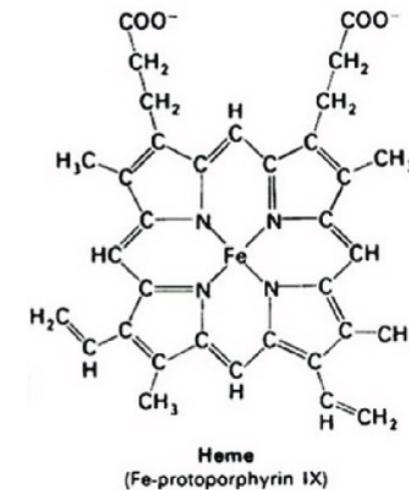
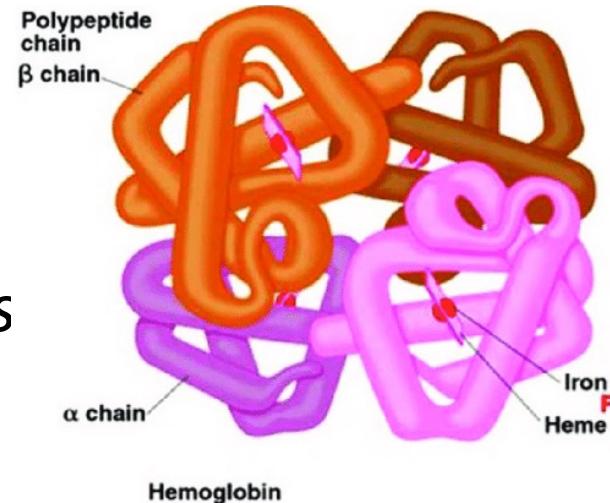
- The study of small molecules (metabolites) produced by a biological system
 - Small usually defined as <1 kda
 - Mass of glucose ~ 0.18 kda
 - Mass of heme ~ 0.62 kda
 - Mass of hemoglobin alpha ~ 16 kda
- Sometimes distinguished from *metabonomics*
 - Includes other molecules in the system, e.g. food
 - Term more common in NMR literature





What is metabolomics?

- The study of small molecules (metabolites) produced by a biological system
 - Small usually defined as <1 kda
 - Mass of glucose ~ 0.18 kda
 - Mass of heme ~ 0.62 kda
 - Mass of hemoglobin alpha ~ 16 kda
- Sometimes distinguished from *metabonomics*
 - Includes other molecules in the system, e.g. food
 - Term more common in NMR literature





Metabolomics philosophies

- **Targeted metabolomics**

- Quantifies 10s of well characterized metabolites from a sample
- Provides absolute quantification through comparisons with chemical standards
- “Usual suspects”: central carbon metabolism, amino acids, nucleotides

- **Untargeted metabolomics**

- Attempts to enumerate 1,000s of unique metabolites from a sample
 - Only characterizing what you can, and with mixed accuracy
- Usually requires multiple screens targeting different chemistries
- Provides (mostly) relative abundances



Metabolomics technologies

- **Nuclear Magnetic Resonance (NMR) spectroscopy**

- Uses environment-specific atomic resonance (e.g. ^1H) patterns to identify compounds
- Can also be used to work out molecular structures (from coupling)
- *Pros:* Absolute, reproducible quantification; minimal sample prep; sample unperturbed
- *Cons:* Low sensitivity (# of molecules, min. abundance); machines very expensive

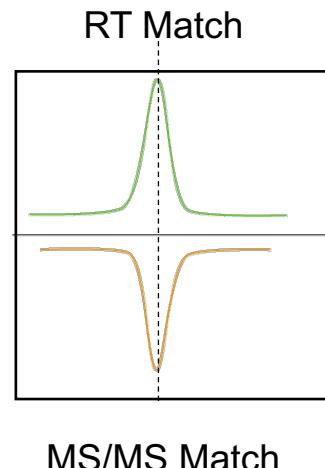
- **Mass spectrometry (MS)**

- Directly measures mass/charge ratio of a metabolite
- Can also be used to work out molecular structure (from fragmentation pattern)
- *Pros:* Very sensitive and flexible; can differentiate 1,000s of metabolites
- *Cons:* Relative quantification only; sample destroyed; complex preps

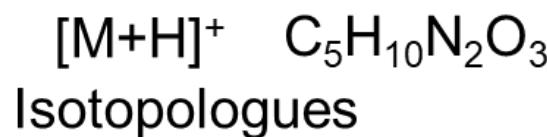


Metabolite identification

- Work out mass (or more properly mass-to-charge ratio, m/z) and thus potential formulas for parent compound
- Match against metabolite databases (e.g. HMDB, Metlin)
 - Likely to be lots of hits
- Match against compound libraries
 - Internal RT values, known/theoretical MS² fragmentation



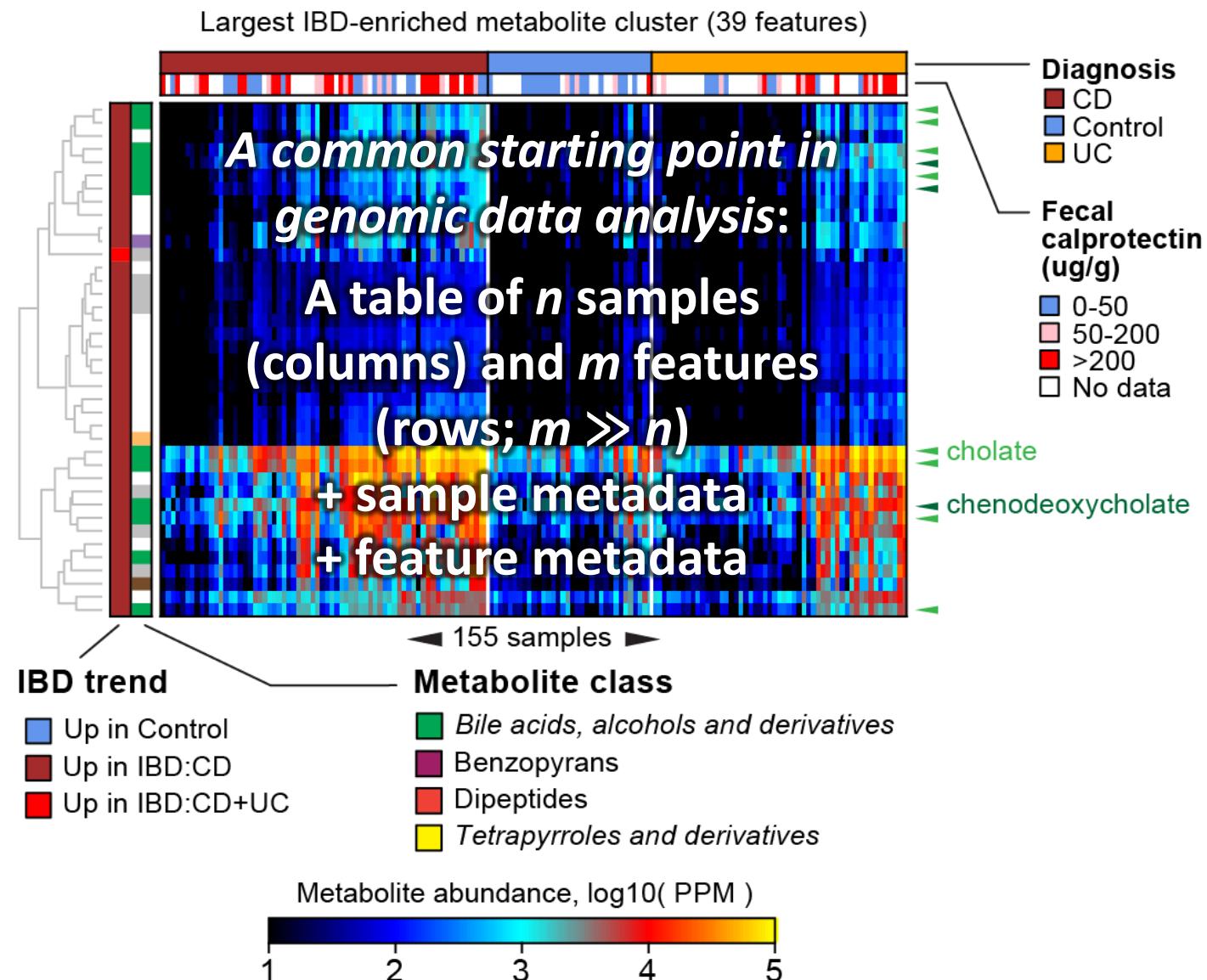
At a given retention time:





What to do with your sample-by-metabolite table?

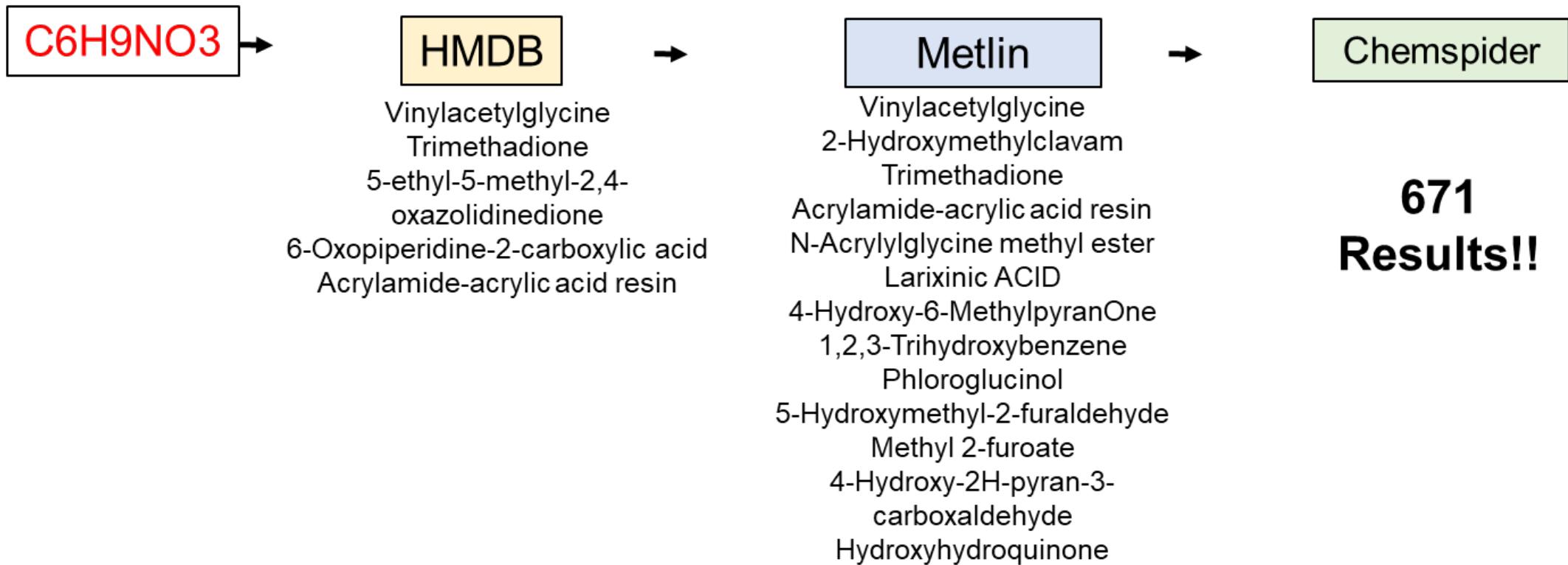
- Quality control
 - Normalization, outlier removal
- Supervised analyses
 - Phenotype association/prediction
 - Enrichment analyses
- Unsupervised analyses
 - Clustering samples
 - Clustering metabolites
- Data integration
 - Associating metabolites with other ‘omics measurements
 - Do characterized metabolites associate with uncharacterized genes?
 - Do uncharacterized metabolites associate with characterized genes?





Metabolite identification

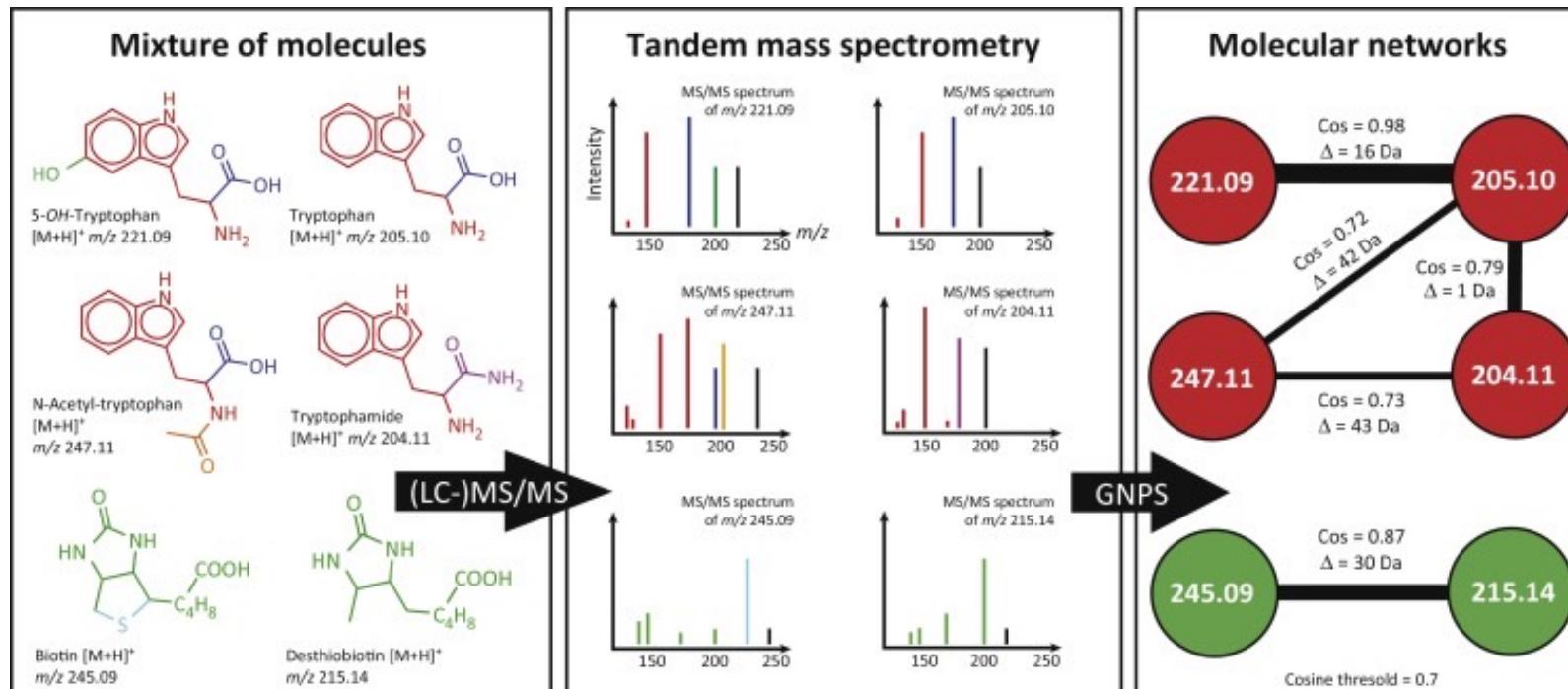
- A formula alone can be very ambiguous (even worse if allowing ± 1 da)





Metabolite identification

- Analyze networks of similar fragmentation spectra (MS/MS only)
 - “Molecular Networking” by Pieter Dorrestein / GNPS
- Transfer knowledge from knowns to unknowns (guilt-by-association)

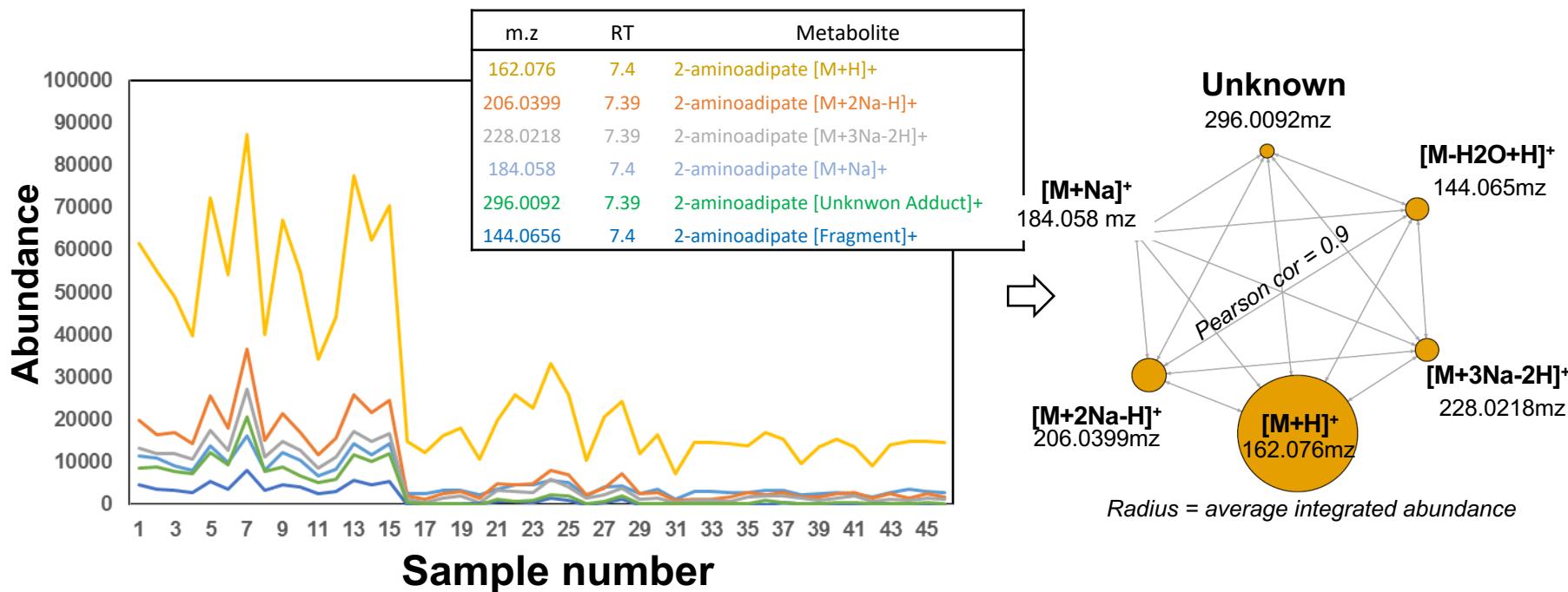


Trends in Pharmacological Sciences



Metabolite identification

- Analyze cross-sample covariation
 - Identify groups of related metabolites (similar to co-expression)
 - Transfer annotations from known to unknown metabolites (also guilt-by-association)



Brown, et al (2011) Bioinformatics 27: 1108–1112

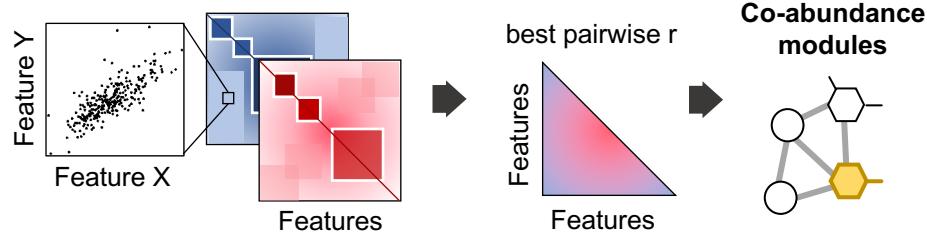
Parent ion
+ isotopes
+ adducts
+ common fragments
(all in same RT group)

MACARRoN: Metabolome Analysis and Combined Annotation Ranks for pRediction of Novel bioactives

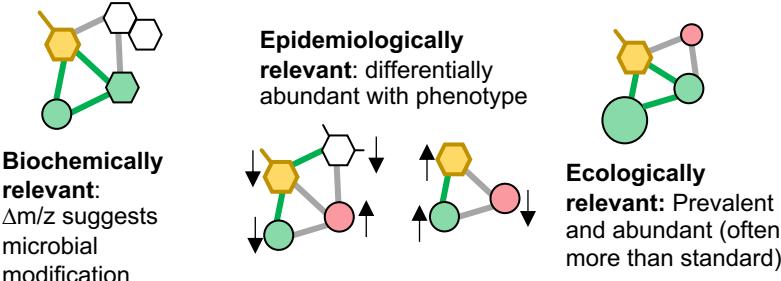


Amrish
Bhosle

Guilt-by-association to transfer potential annotations



Signs of potential bioactivity

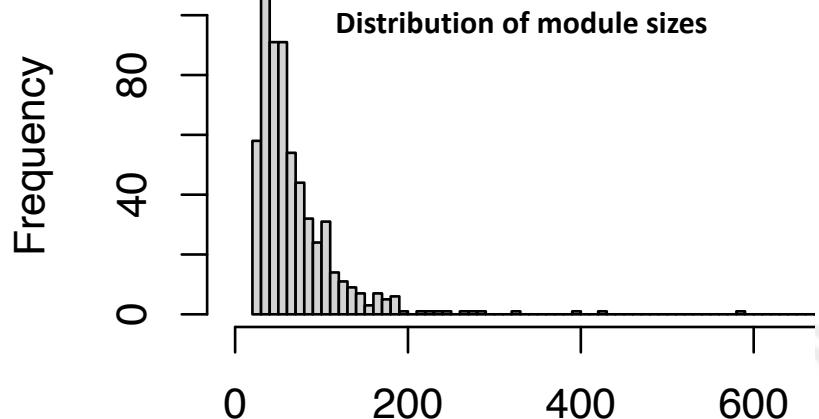


546 HMP2 fecal metabolomes

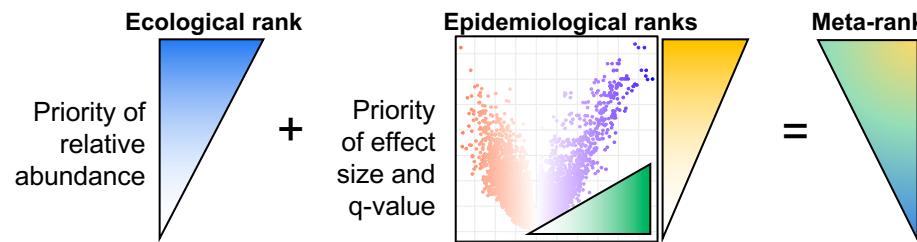
(non-IBD + CD + UC)

- 44,757 features
- **466 annotated**

43,498 features in 606 modules

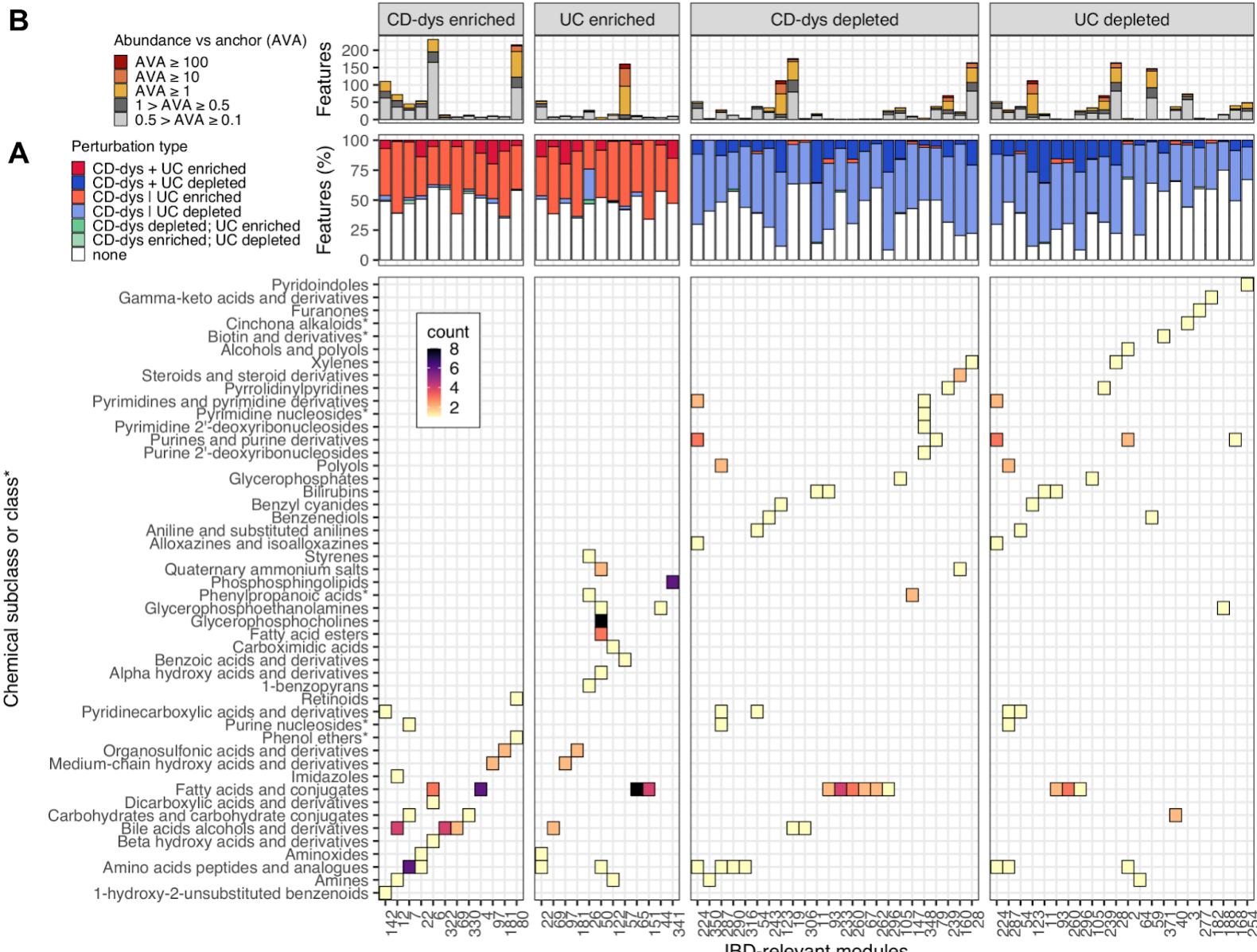


Integrated prioritization of potential bioactives

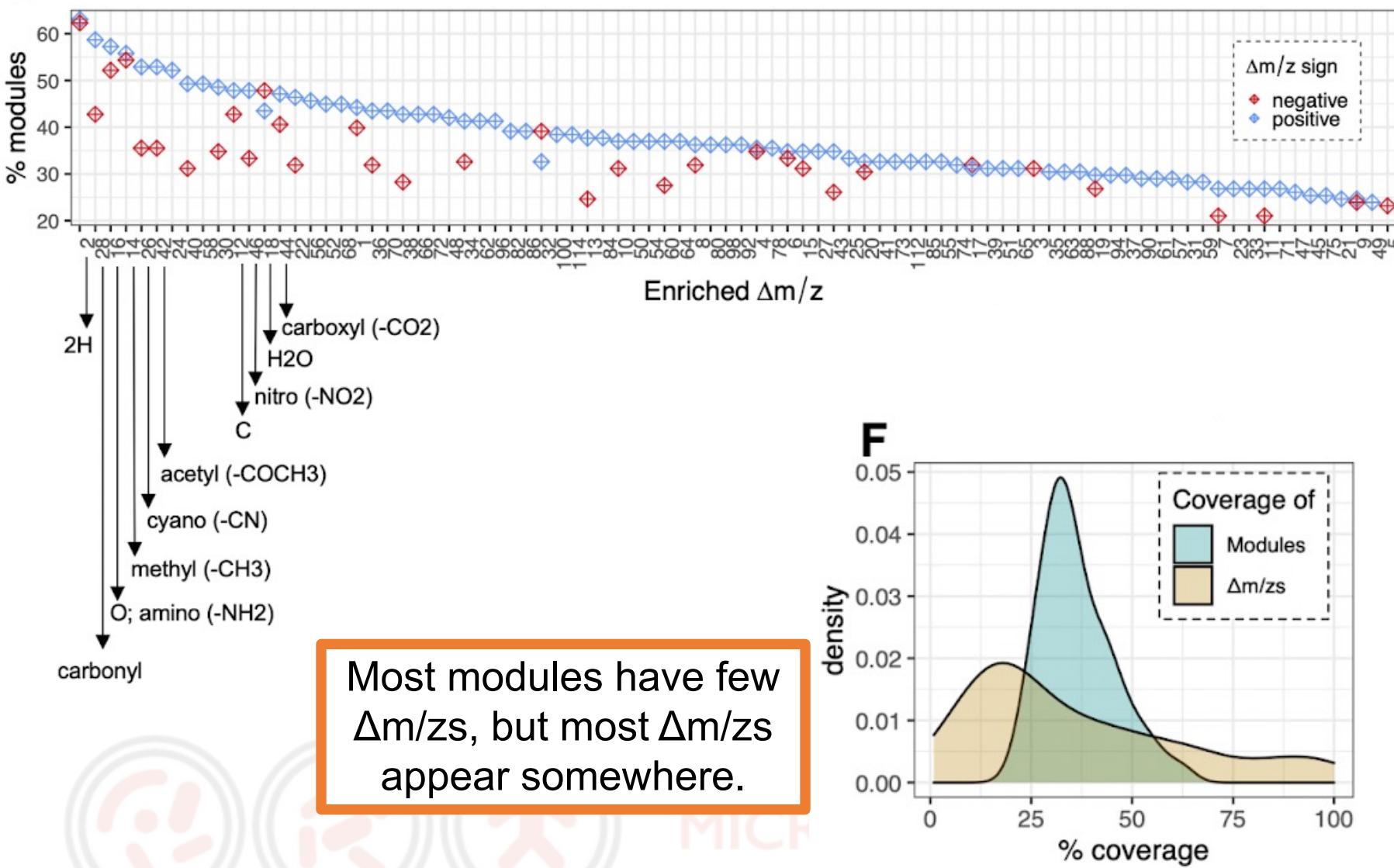


152 modules with at least one annotated metabolite →
13,633 features associated with at-least one annotated metabolite

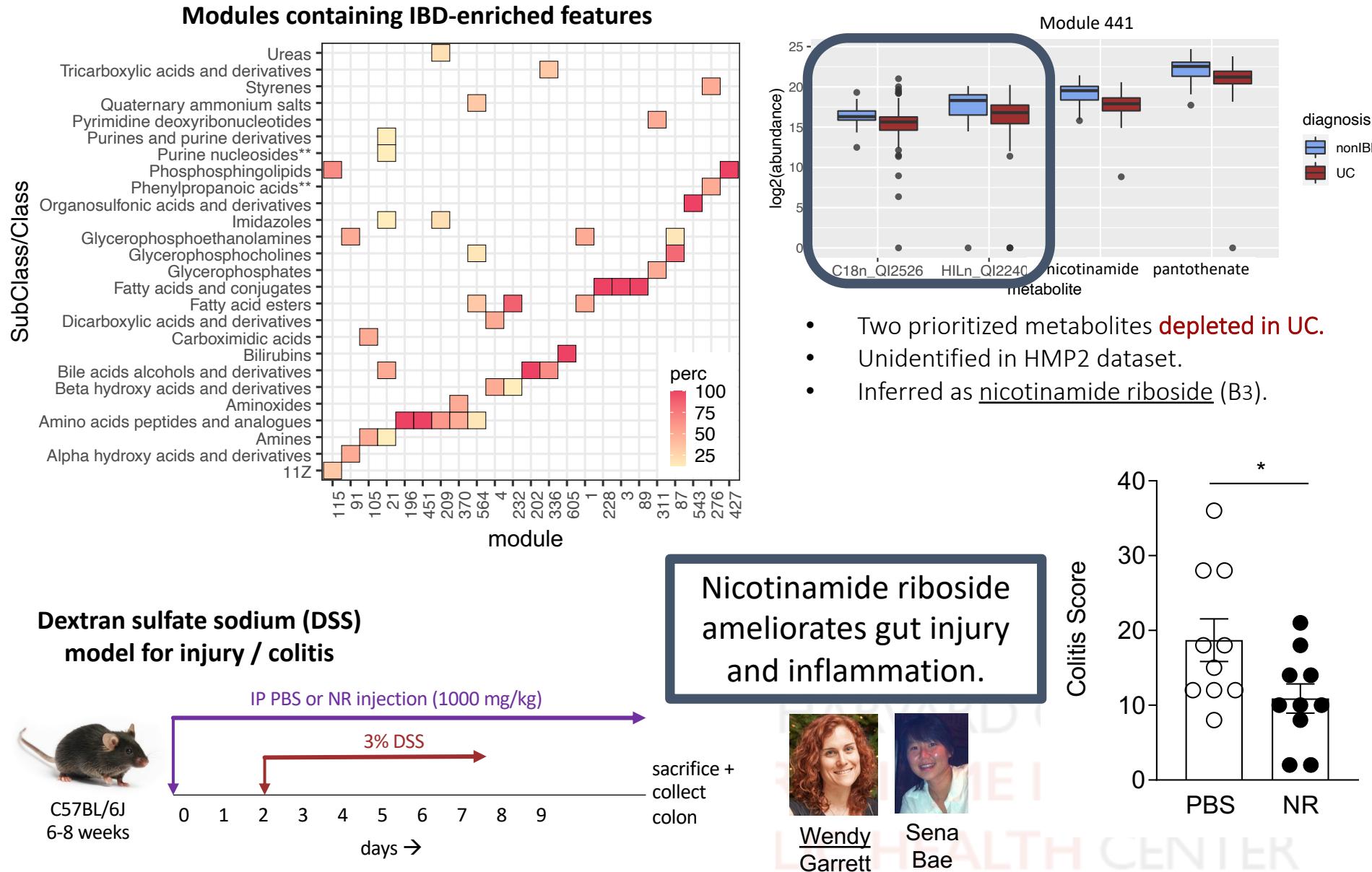
A large subset of modules are IBD-relevant and characterizable



Putative modifications are characteristic of specific parent compounds



Bioactivity predictions are consistent and (at least one is) validated *in vivo*





Summary

- Metatranscriptomics (MTX)
 - Functional profiling of MTX data
 - Relationships between functional potential and activity
 - Detecting under- and over-expression
 - Statistical models for MTX analysis
- Metabolomics (MBX)
 - Sequencing is easy... everything else is hard
 - Metabolomics approaches
 - Identifying compounds
 - Bioactivity prediction with MACARRoN



(One way of thinking about)

Strains in microbial communities



Summary

- What is a microbial strain?
- Why is considering strains important?
- Gene-level strain profiling with PanPhlAn
- SNV-level strain profiling with StrainPhlAn
- Community strain statistics with ANPAN



Two big questions of shotgun meta'omics

Who is there?

Taxonomic profiling of (e.g.) phyla, species & strains



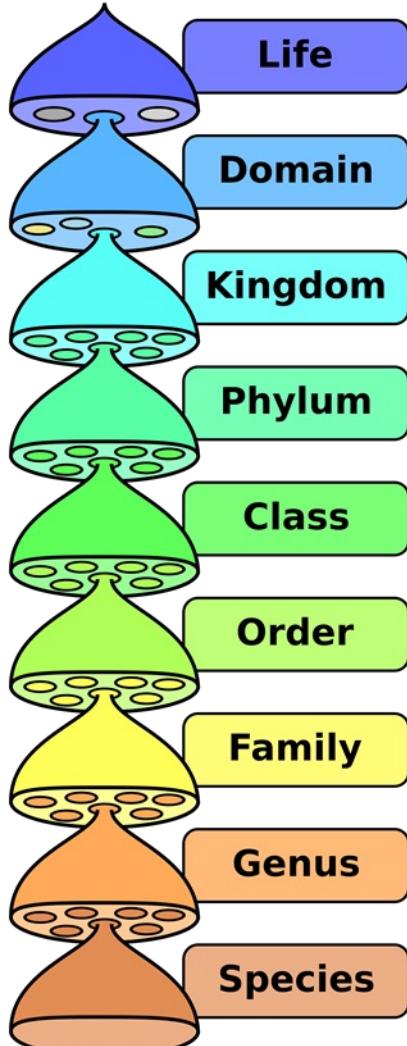
What are they doing?

Functional profiling of (e.g.) enzymes & pathways

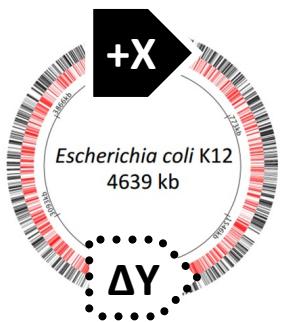




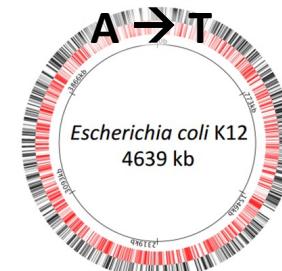
What is a microbial strain?



- A microbial strain is a taxonomic group below the species level
 - Defining species is not trivial, though systematic definitions exist
 - 95% average nucleotide identity; 97-99% 16S rRNA identity
- Strain = population of clonal cells
 - OK for culture applications, but less applicable to sequencing
- Strain = specific microbial isolate genome sequence
 - Very precise, but often too specific to be useful
- Strain = collection of “very similar” cells/genomes
 - With “very similar” defined using phylogenetic criteria



*Does adding a gene to a species' genome make a new strain?
How about deleting one?*



*What about changing a single nucleotide?
What about 10 or 100?*



Why care about strains? Consider *E. coli*

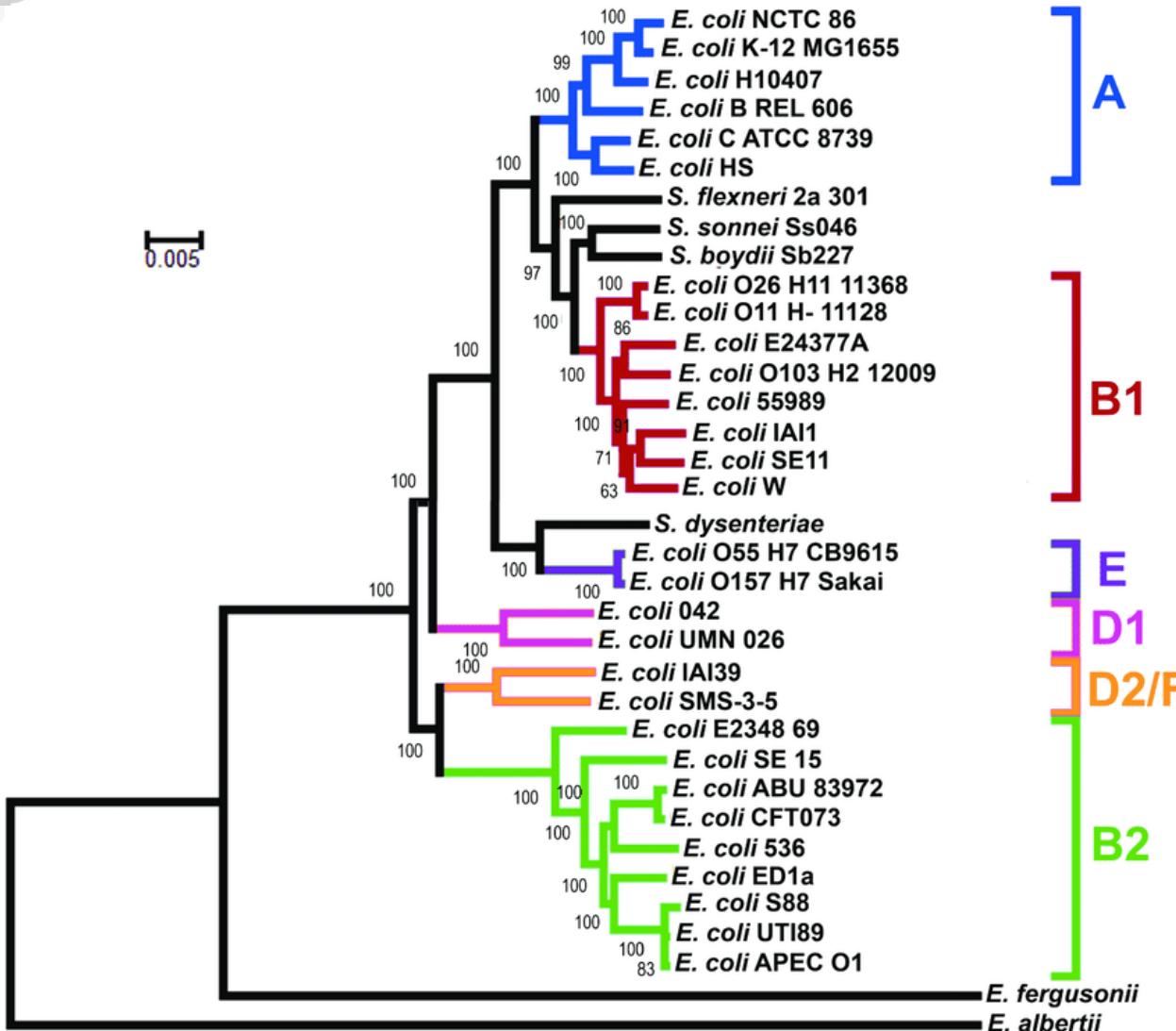
THE MOST ENCOURAGING THING IS THIS: THE STRAIN OF *E. COLI* USUALLY USED FOR CLONING GENES HAS GROWN SO 'DOMESTICATED' DURING ITS YEARS IN THE LAB, THAT IT CAN NO LONGER SURVIVE IN THE HUMAN GUT!!



(Gonick & Whellis 1991)



Why care about strains? Consider *E. coli*



Dunne, Microbial Genomics, 2017

A

B1

E

D1

D2/F

B2

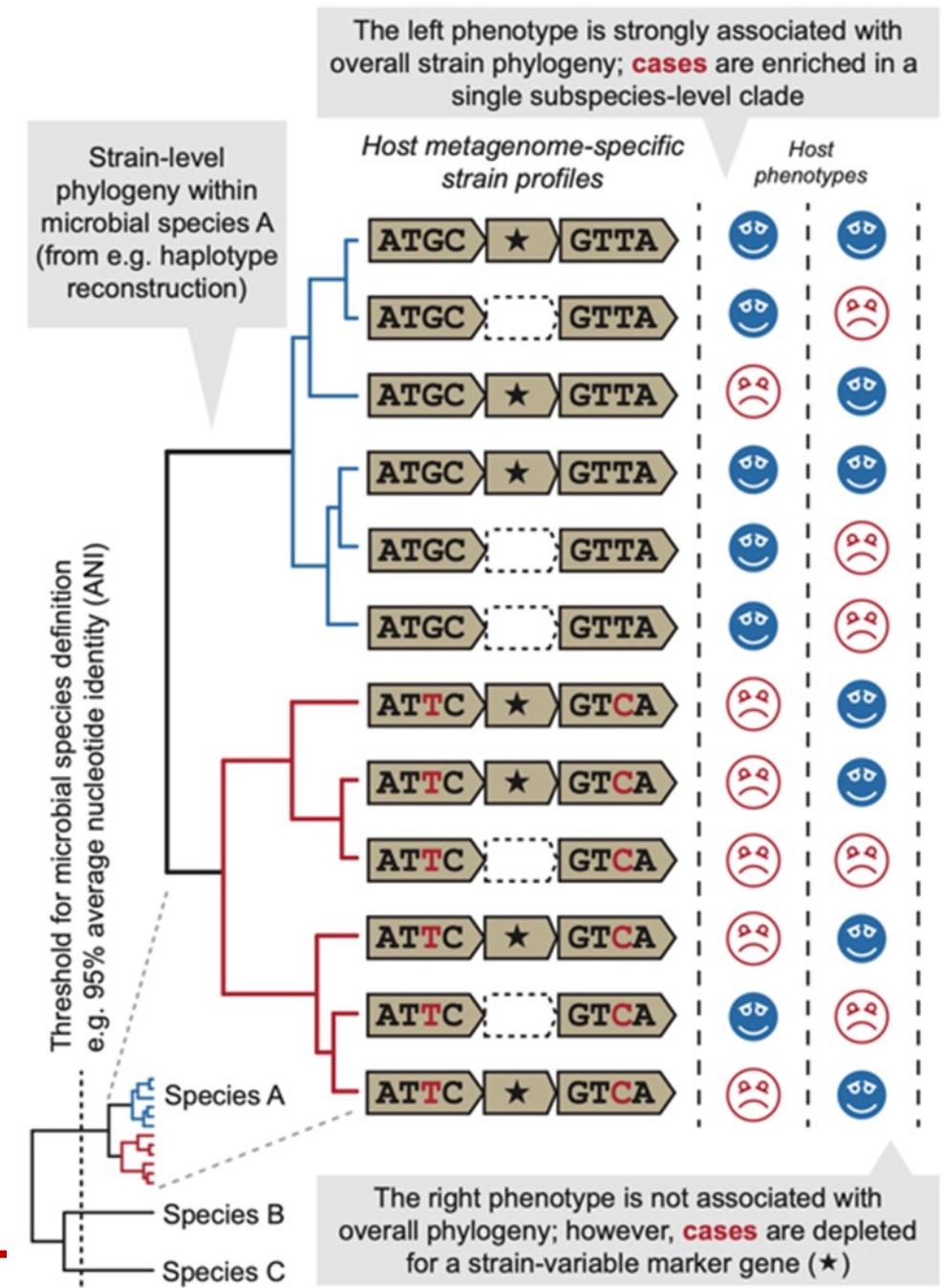
- 1,000s of genomes sequenced
 - Clear sub-species structure
- Some strains are benign
 - Common in the human gut
- Some are acute pathogens
 - e.g. *E. coli* O157:H7
- Some are risk-associated
 - e.g. pks+ *E. coli* in colon cancer
- Some are probiotic
 - e.g. *E. coli* Nissle 1917

What if we looked at other microbial species in the same detail as *E. coli*?



Why care about strains?

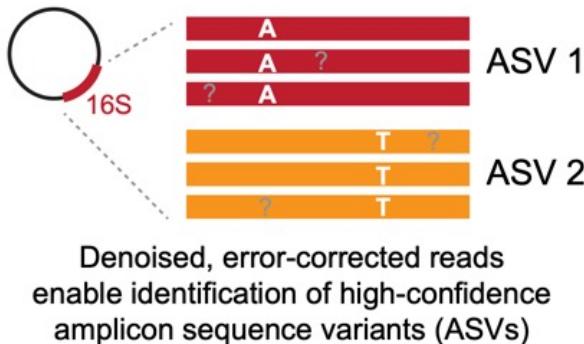
- Strain-level variation can associate with a phenotype or environmental property even when species do not.
- A phenotype ☹ can associate with species phylogeny below the species level, or...
- ...a phenotype ☹ can associate with carriage of a specific gene.
- *(Note that neither of these associations requires defining clear sub-species.)*



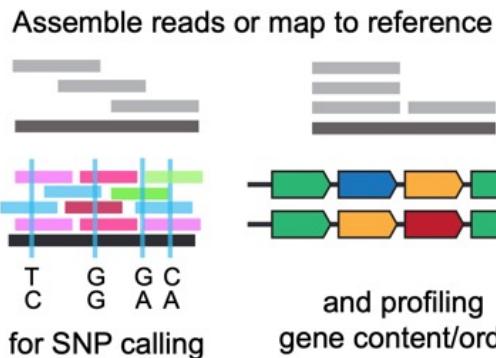


How to study strains in a culture-free context?

A Amplicon sequencing



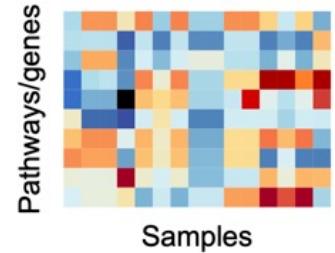
B Metagenomic sequencing



D Single-cell sequencing

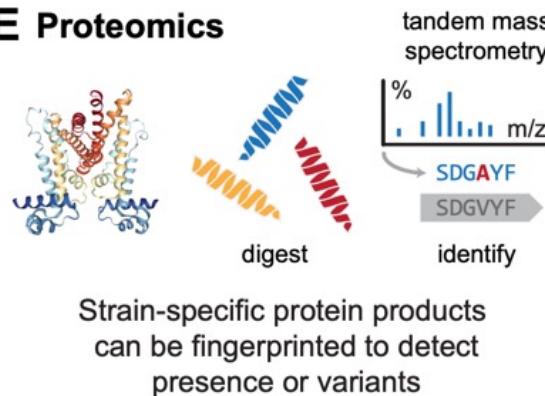


C Metatranscriptomics

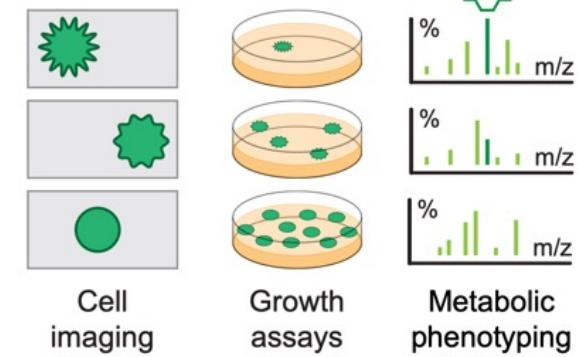


Quantify microbial transcripts by mapping or *de novo* assembly; identify strain-specific differences in gene expression.

E Proteomics

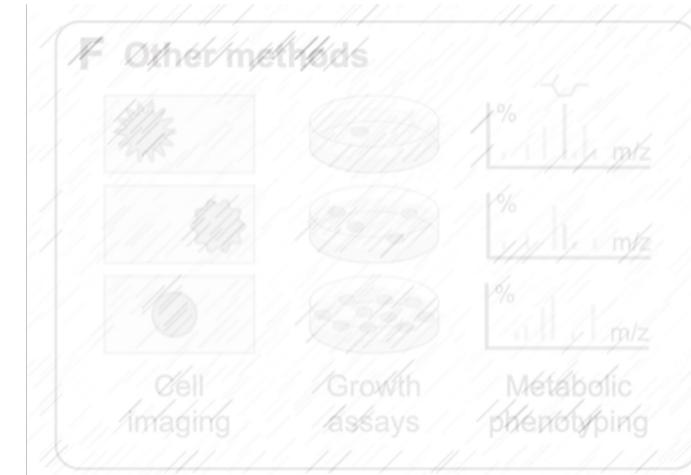
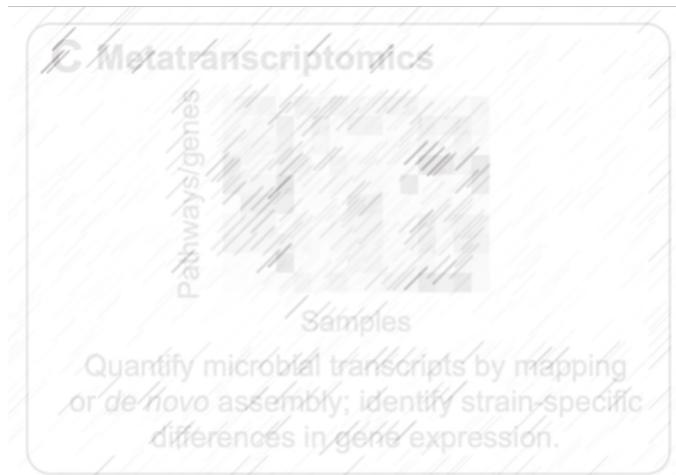
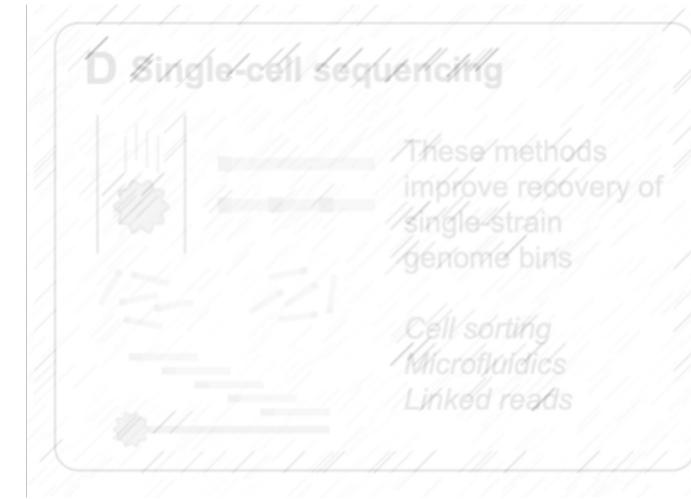
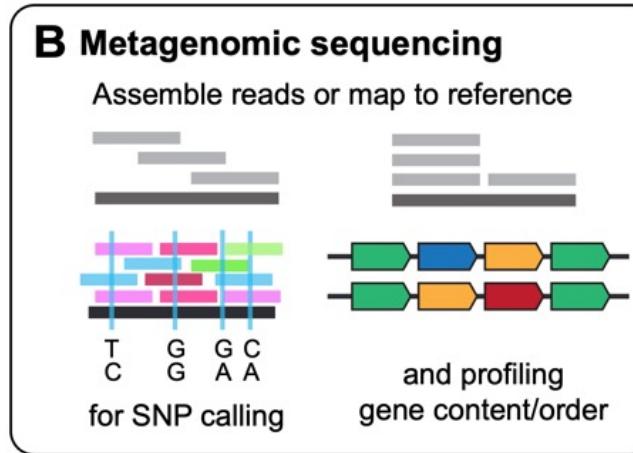


F Other methods





How to study strains in a culture-free context?





Reference-based strain identification

- Four broad mechanisms for reference level strain profiling
 - Identification of reference genotypes in a community (e.g. [PathoScope](#); [Sigma](#))
 - Identification of the dominant strain (incl. novel) per species (e.g. [StrainPhlAn](#), [MetaMLST](#), [MetaSNV](#))
 - Identification of more than one strain per species (e.g. [ConStrains](#), [DESMAN](#), [StrainGR](#))
 - Identification of structural variants within a community (e.g. [PanPhlAn](#))

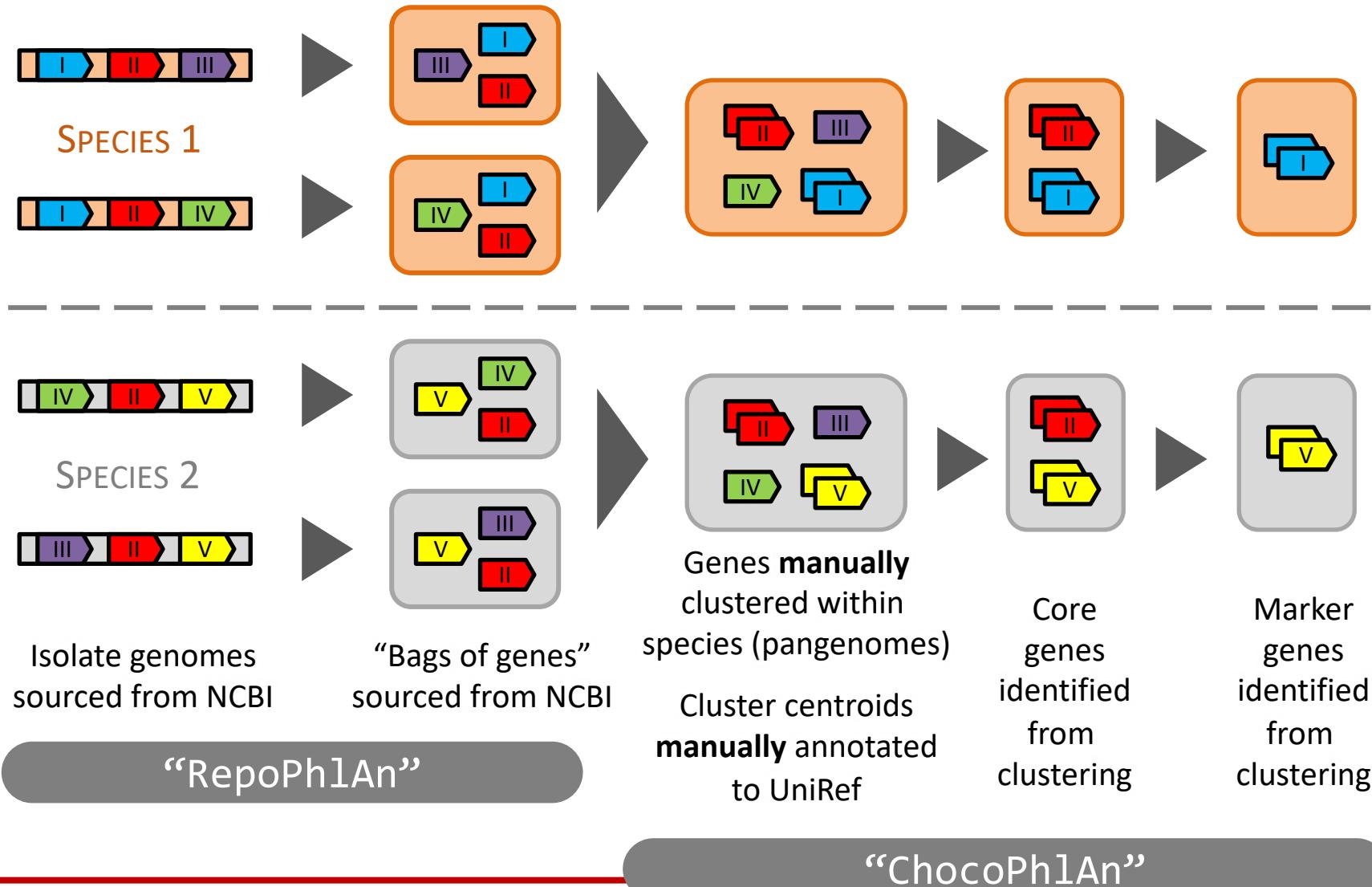




Gene-level strain profiling with PanPhlAn

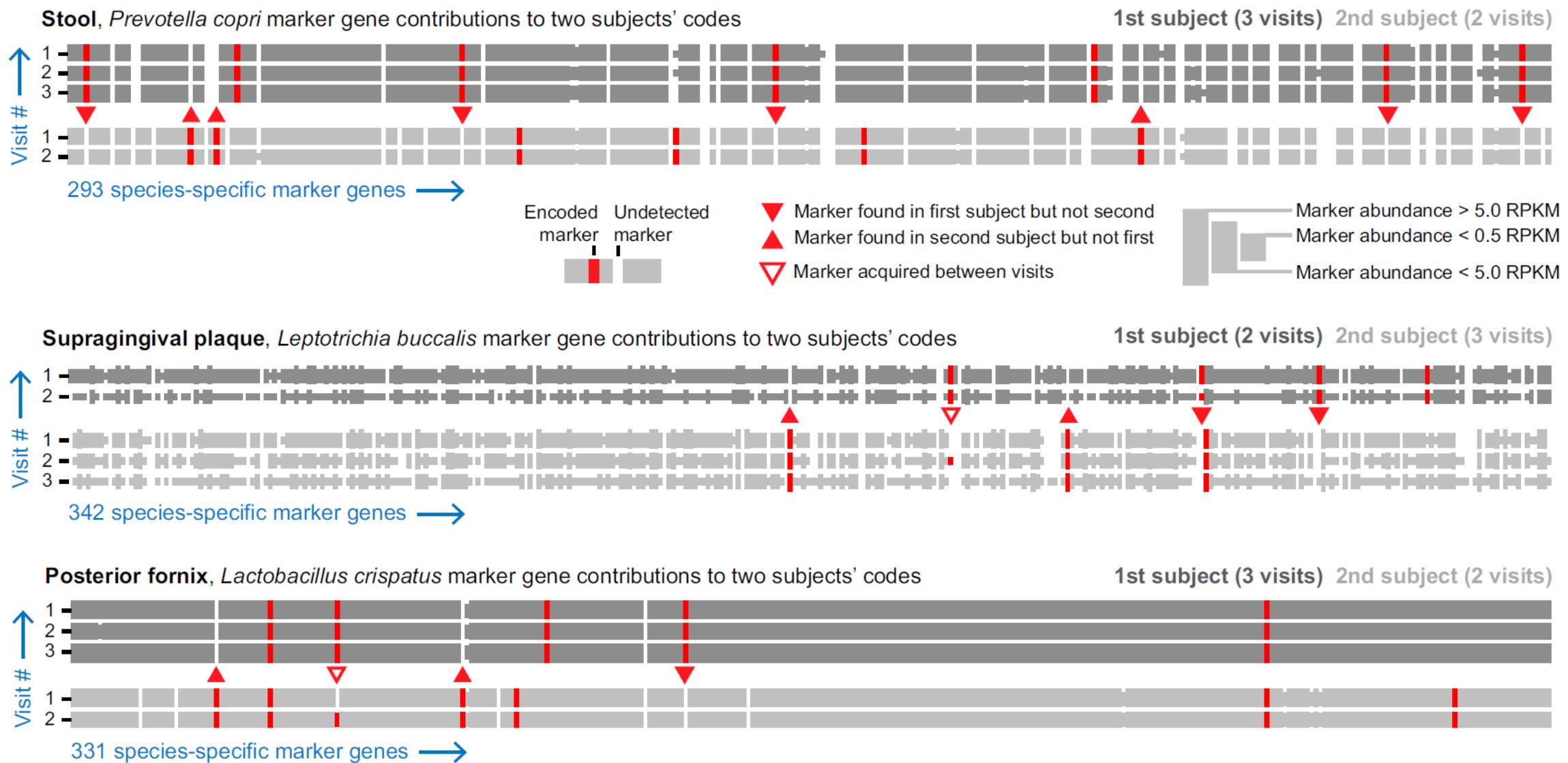


Pangenomes and marker genes as strain profiling tools



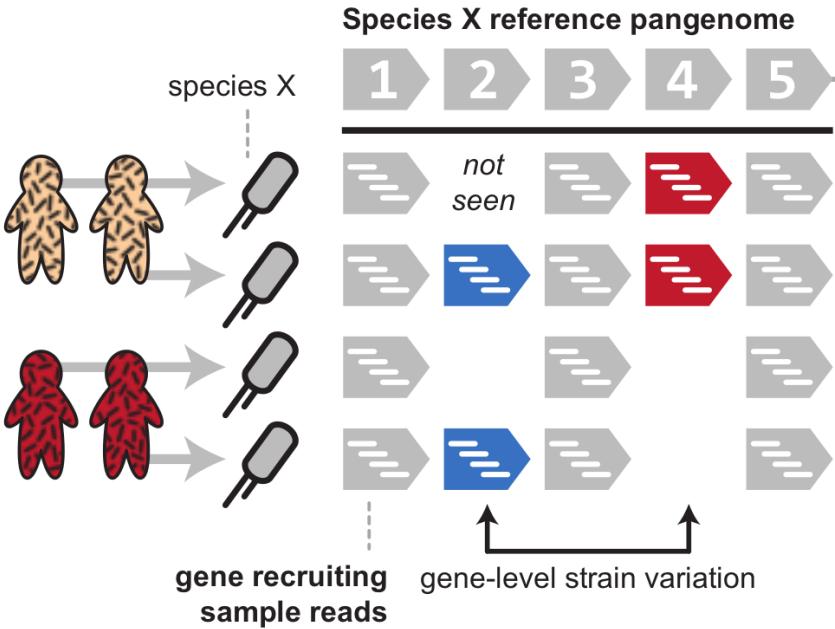
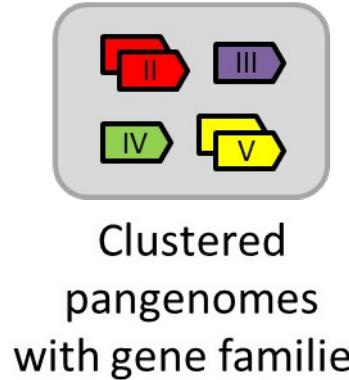
- RepoPhlAn and ChocoPhlAn are database systems for identifying markers genes from isolate genomes (among other tasks).
- We will use other products of these systems in subsequent analysis methods.

Stable, individualized marker loss hints at strains

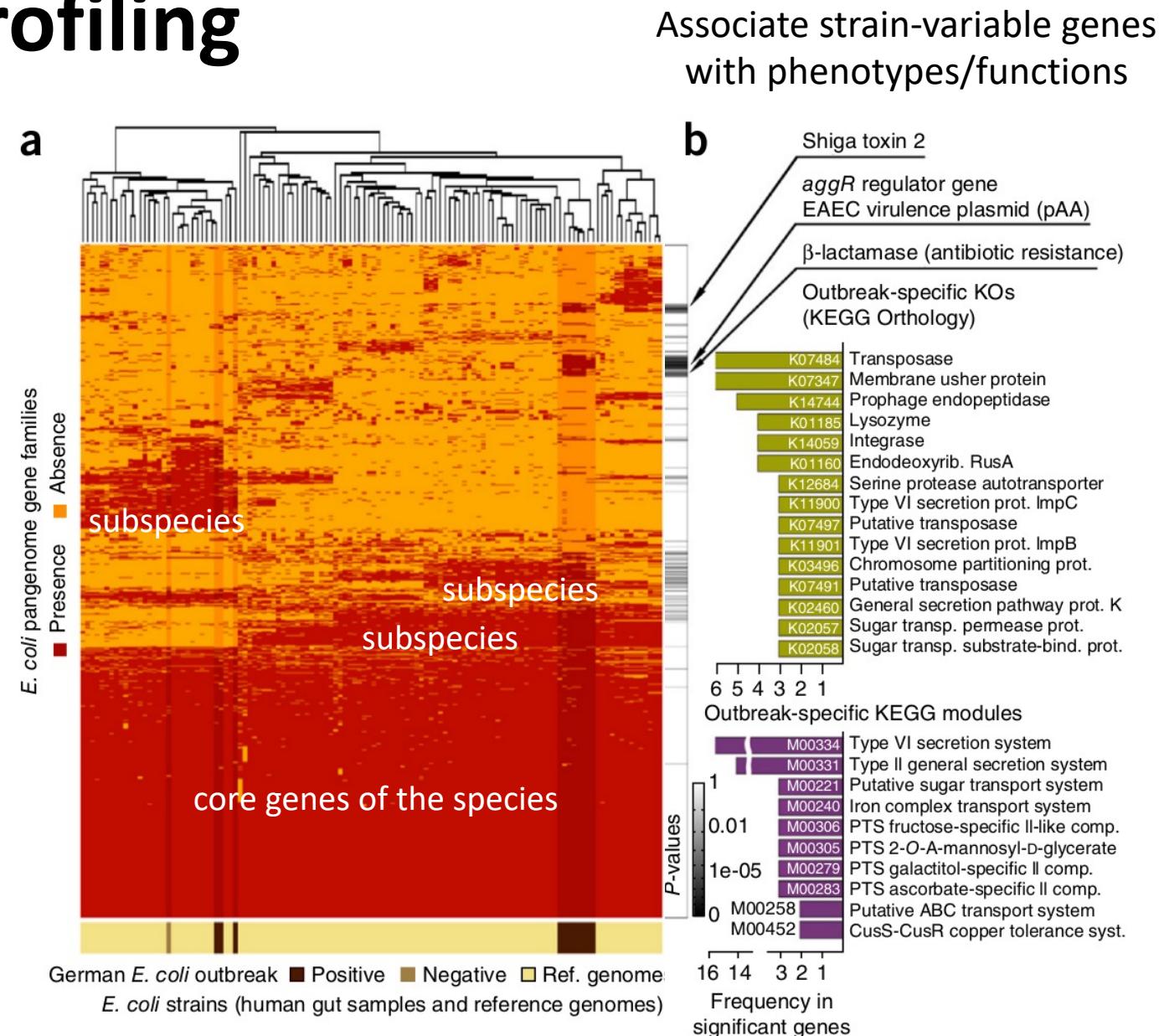




Gene-level strain profiling

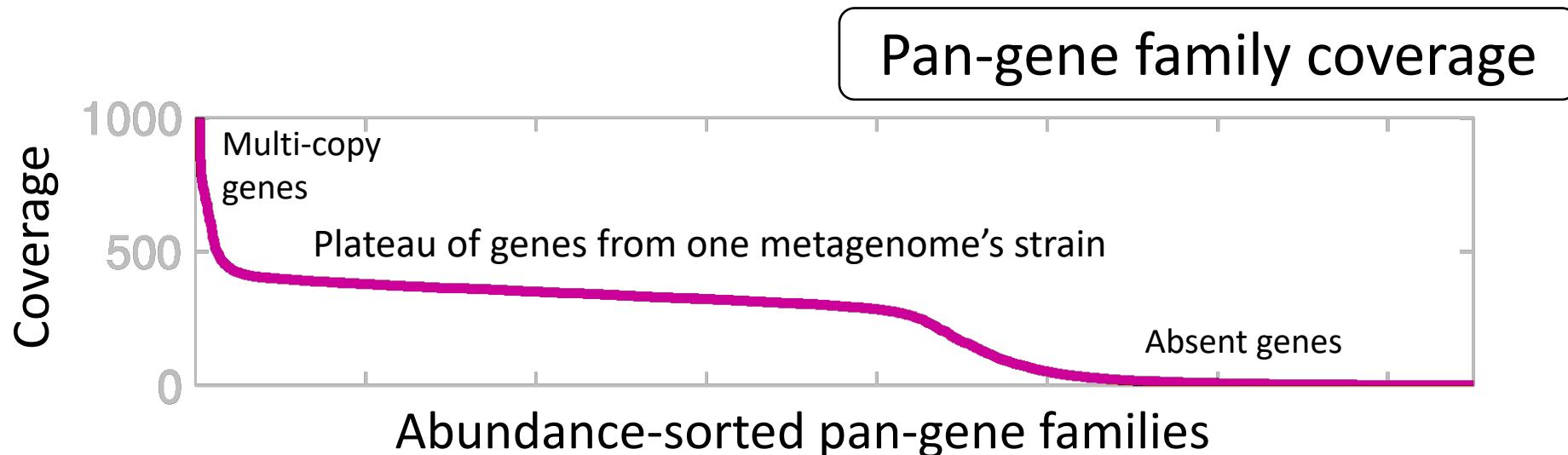
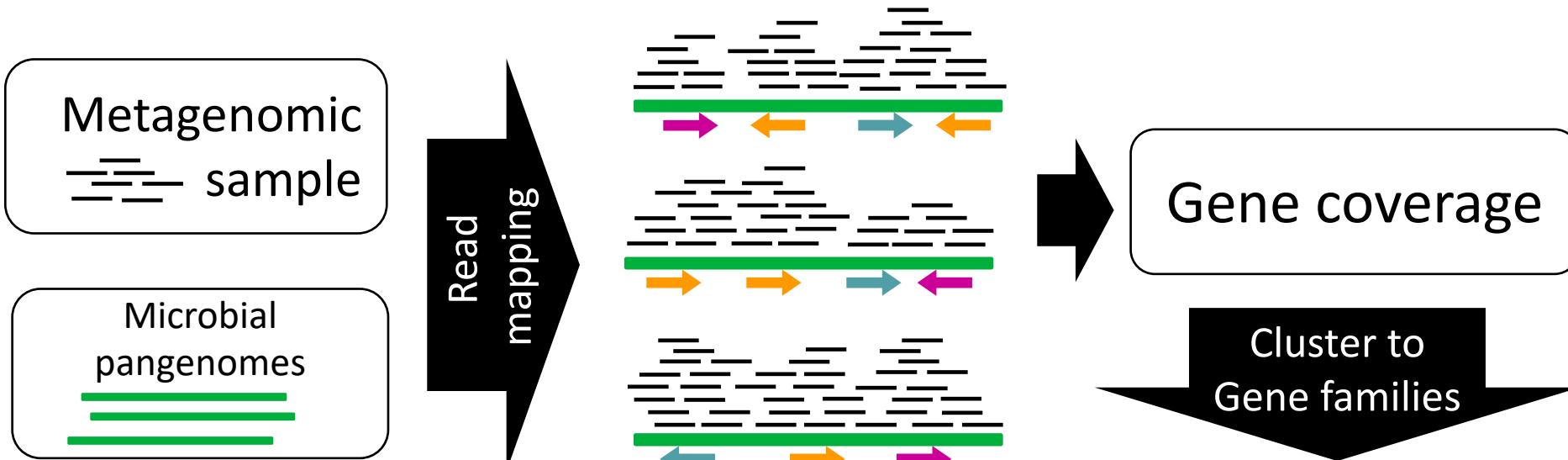


Look for variation in gene presence across samples after mapping reads to species' pangenomes



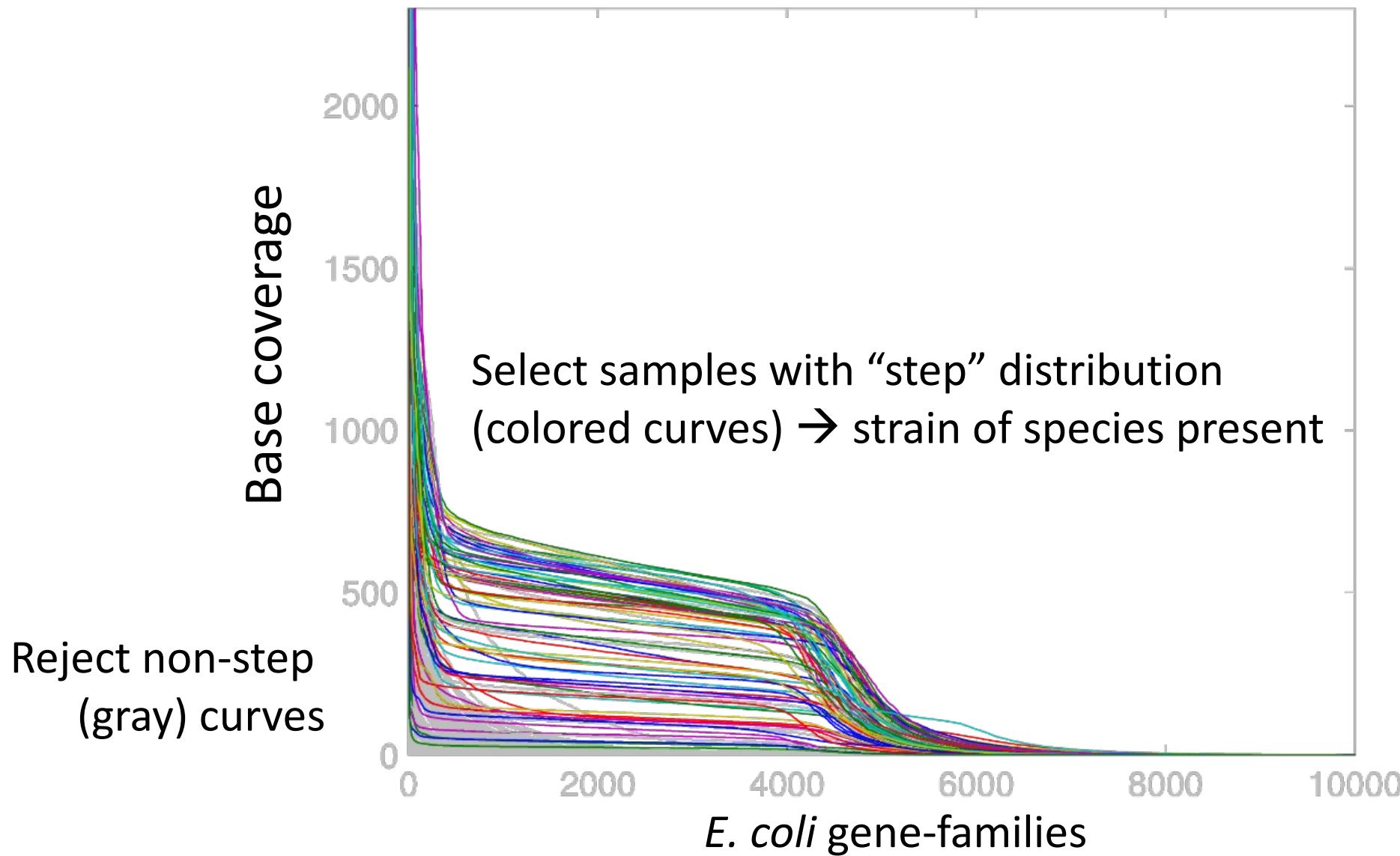


The PanPhlAn approach



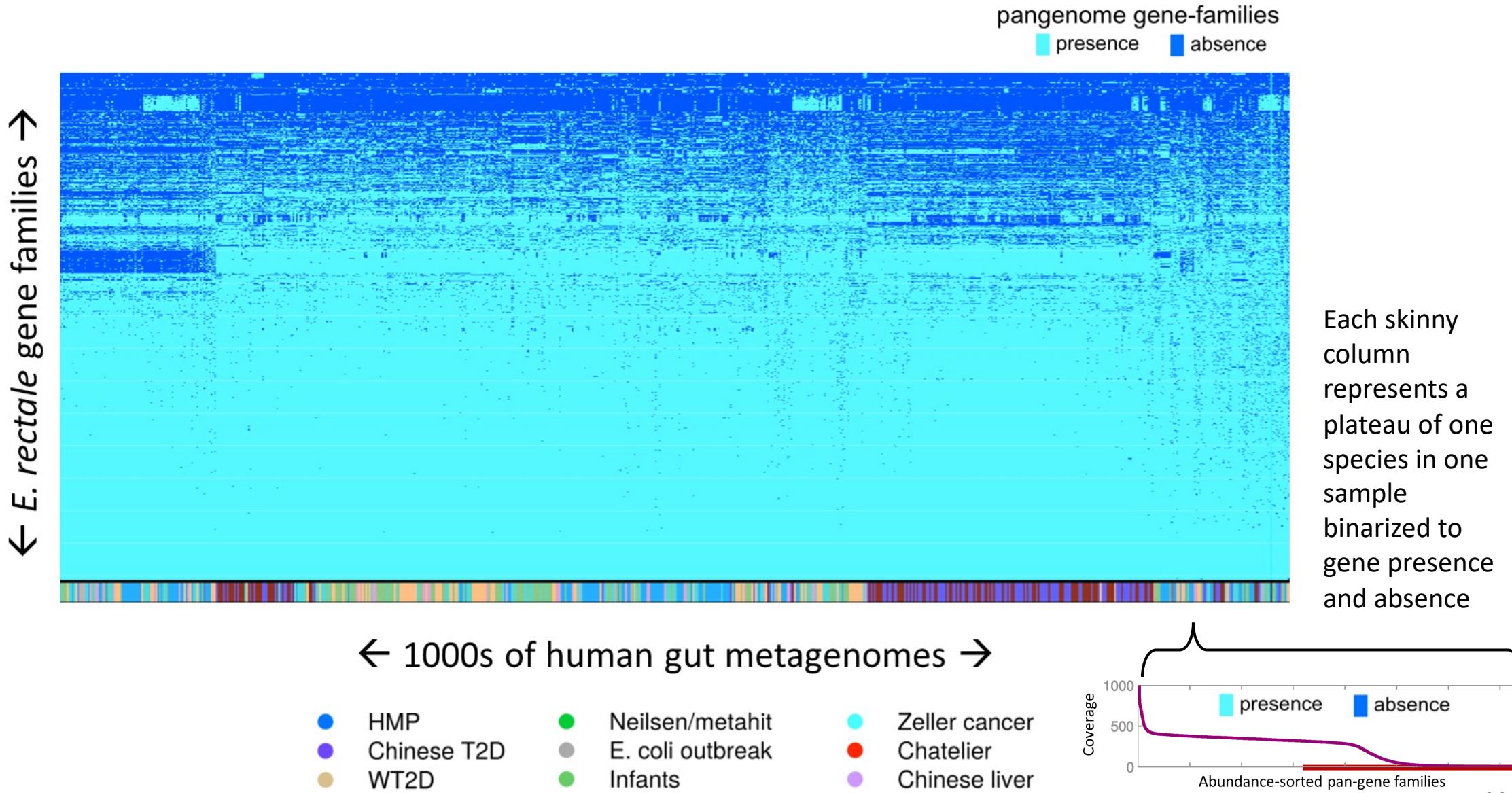


The PanPhlAn approach



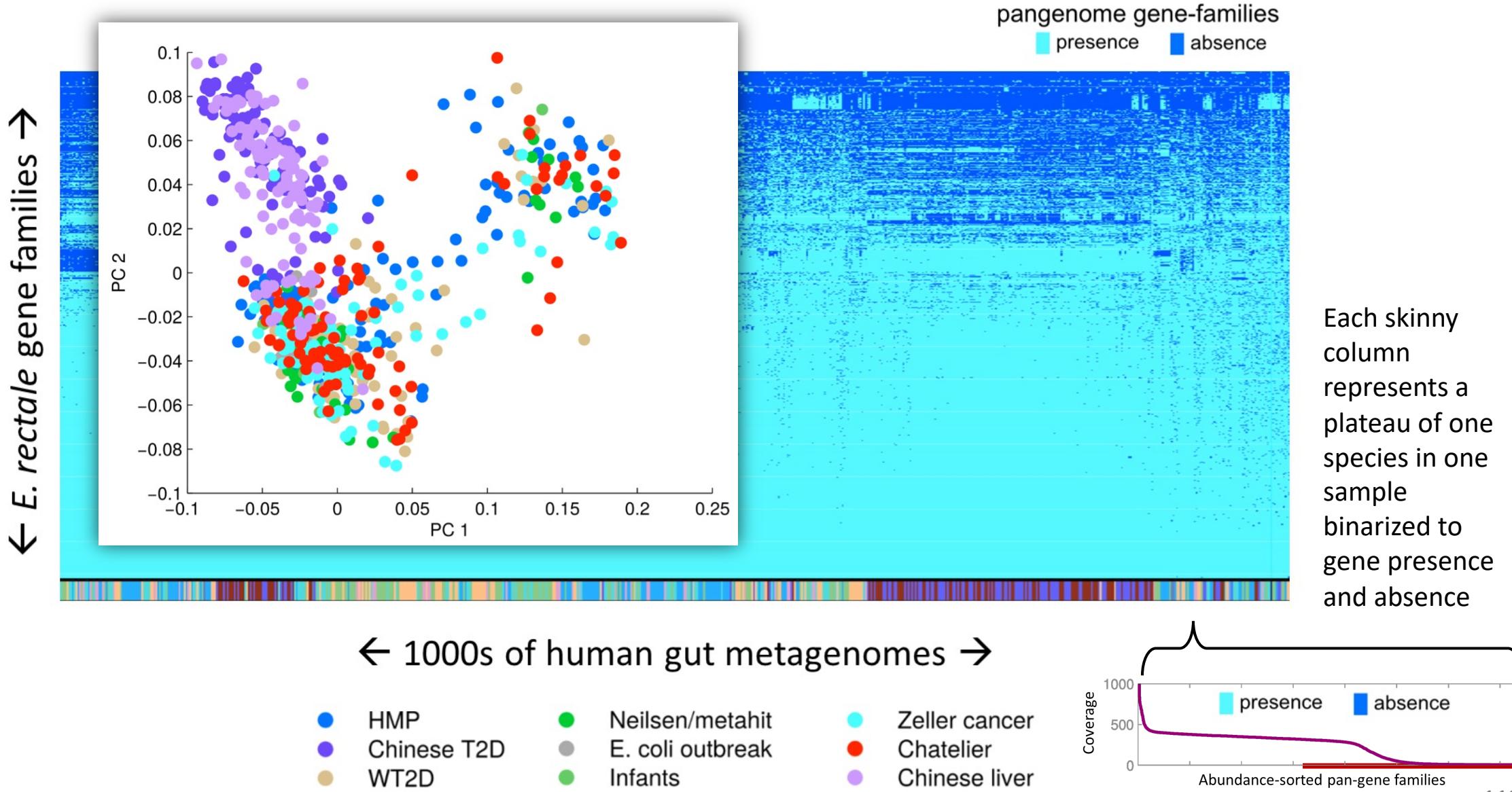


PanPhlAn reveals strain-level population structure of *Eubacterium rectale* from a single reference genome





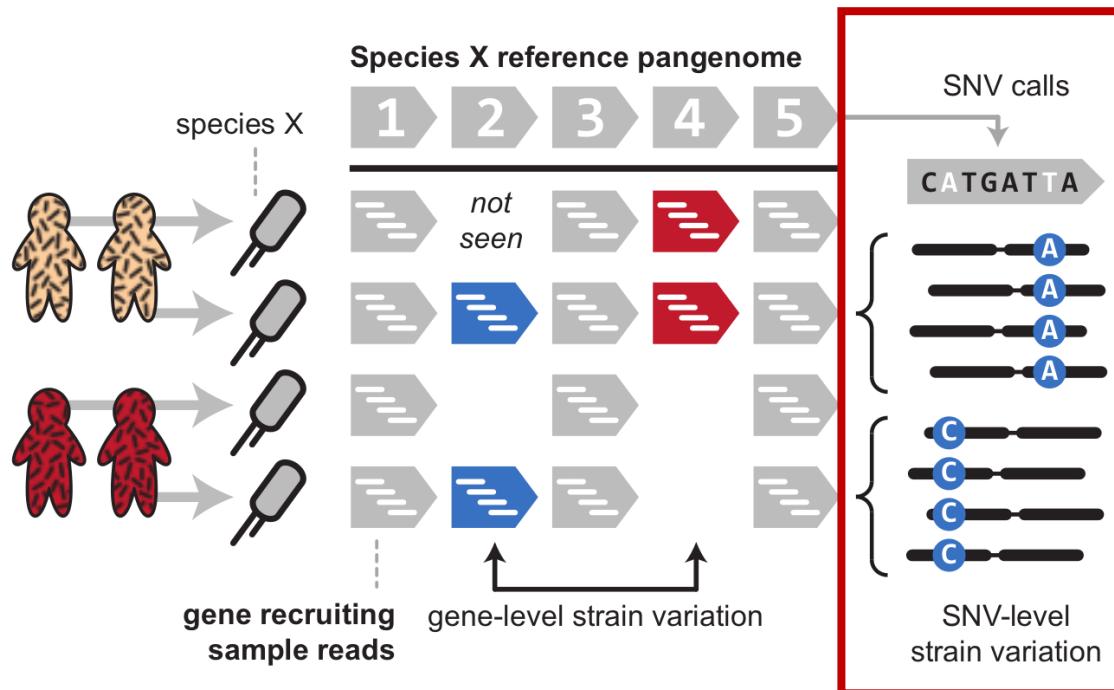
PanPhlAn reveals strain-level population structure of *Eubacterium rectale* from a single reference genome





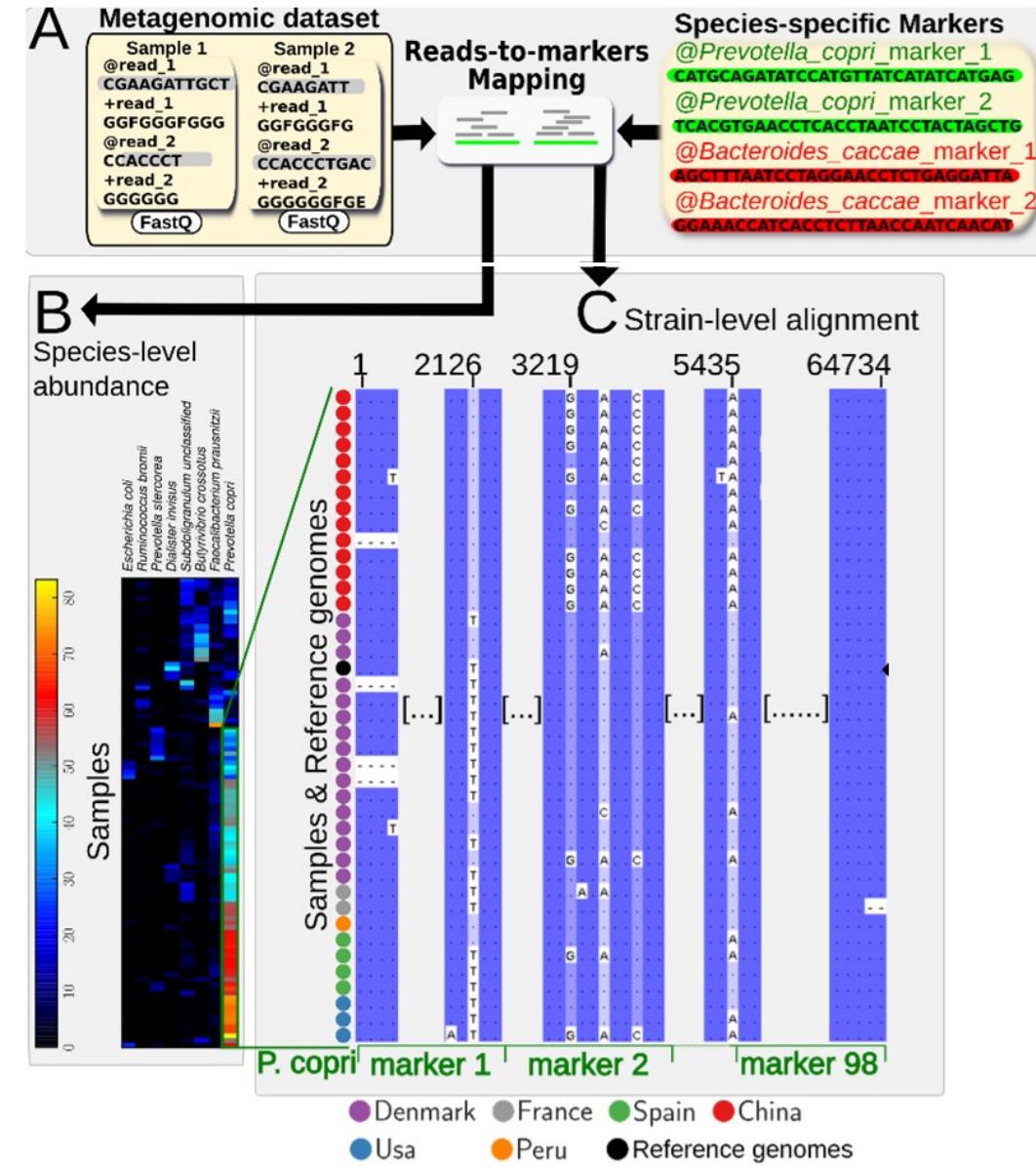
SNV-level strain profiling with StrainPhlAn

Profiling strains at nucleotide-level resolution



- **StrainPhlAn** reconstructs species' dominant per-sample haplotype from SNVs in species-specific marker genes
- For each species, build a strain distance matrix over samples

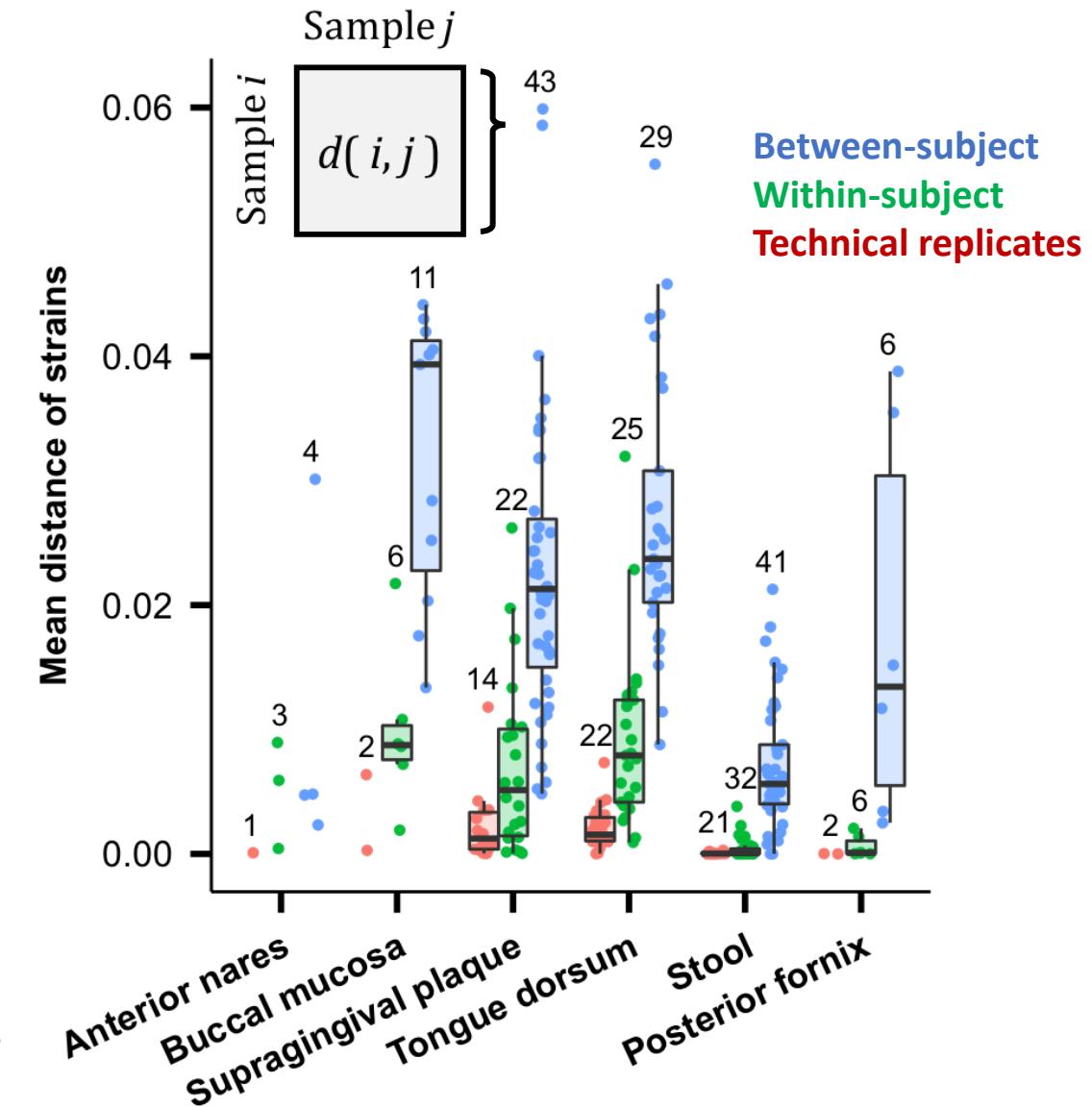
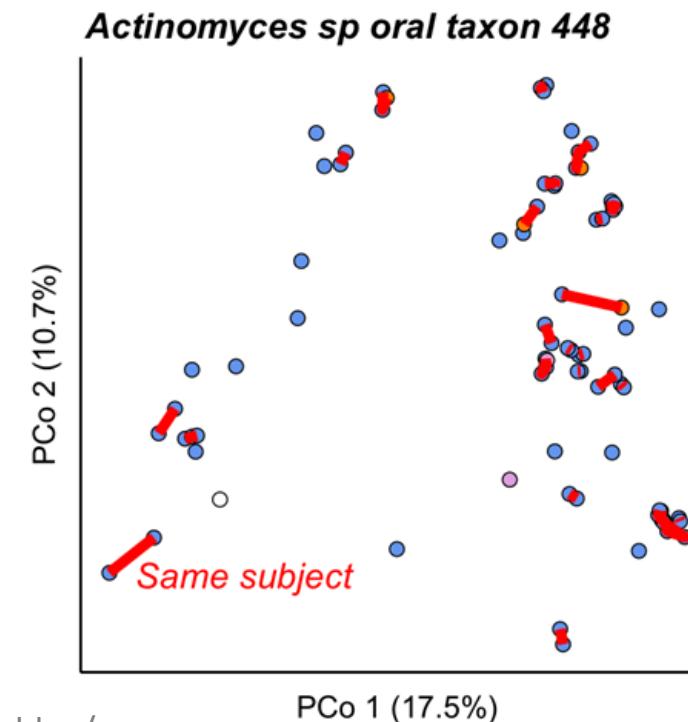
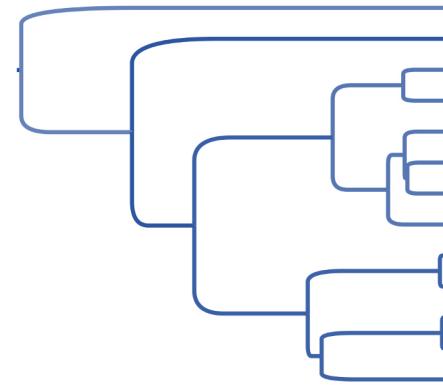
$$\text{Sample } i \quad d(i, j)$$





Strain-level properties of the human microbiome

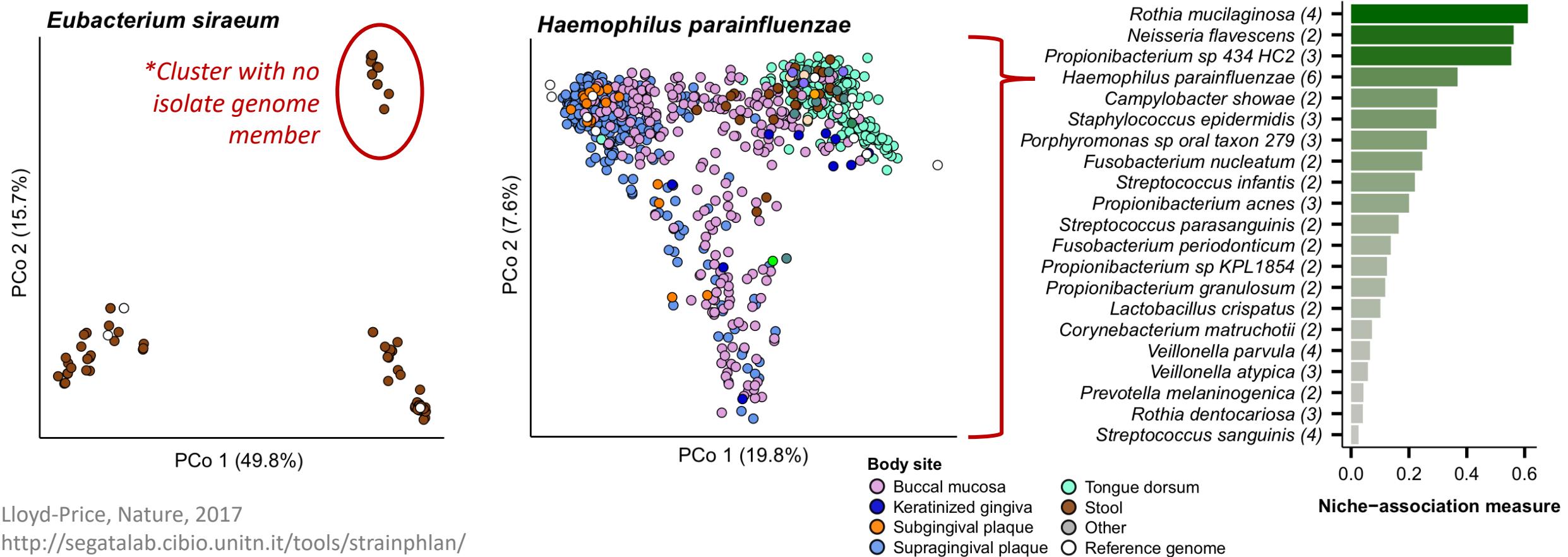
- Less haplotype diversity among stool species compared to oral sites
- $d(\text{between-subject})$
 $\gg d(\text{within-subject}) > d(\text{tech. reps.})$





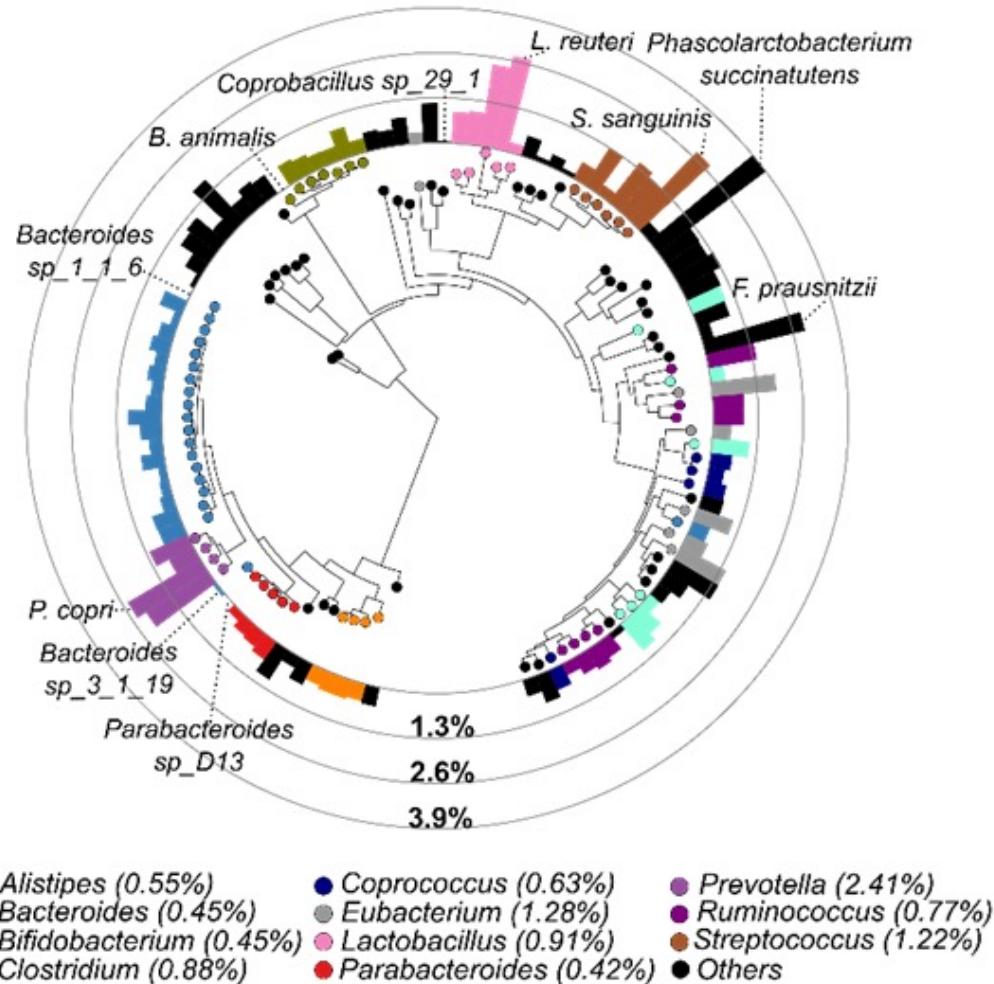
Strain-level properties of the human microbiome

- Subspecies structure manifests as “clusters” of strains
- Structure associated with body site suggests possible niche-adaptation

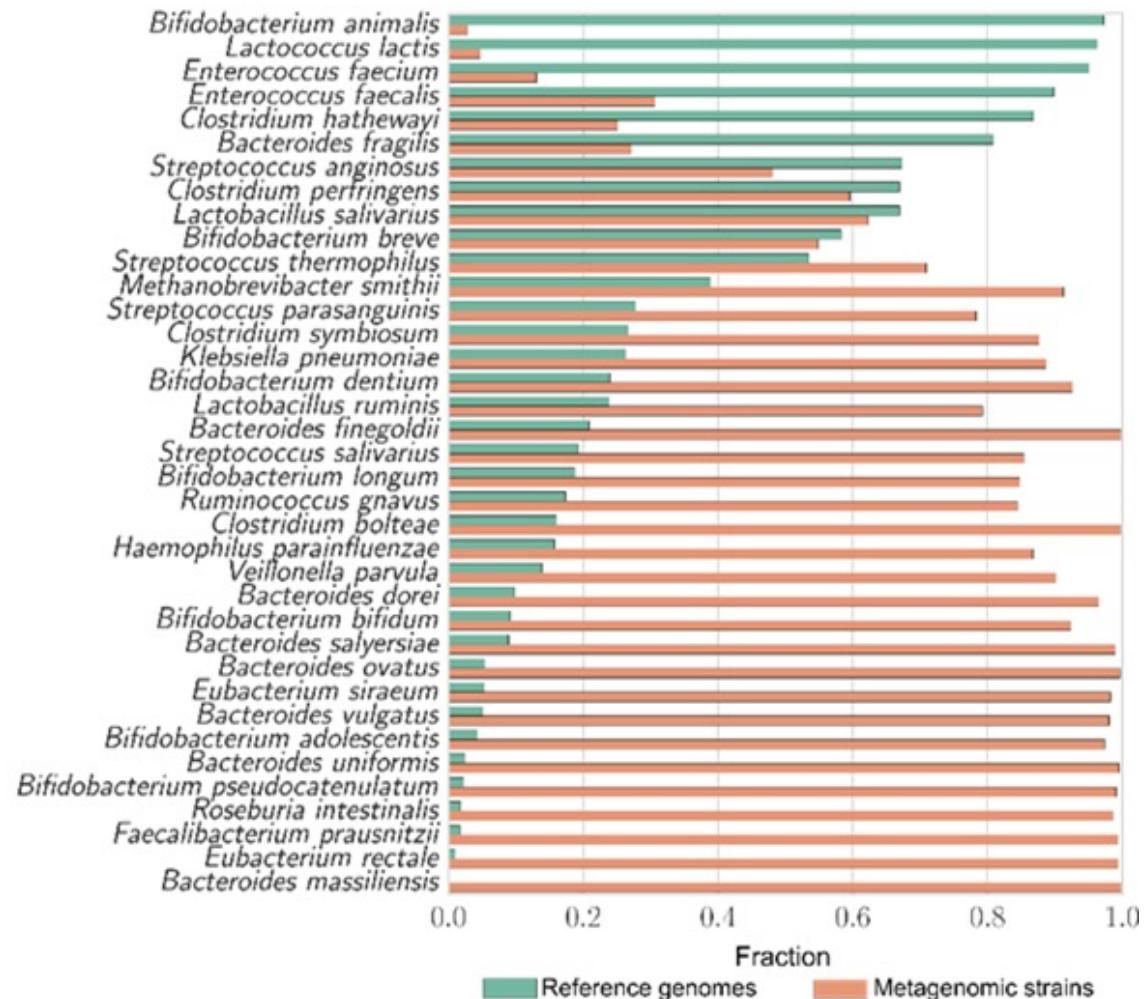




Lots of undersampled species genome diversity



Median divergence from
reference markers



Fractions of intraspecific branch length
for genomes/metagenomic strains



Metatranscriptomics



Strain tracking by gene expression

- Strain identification by analogous differences in metatranscriptomic gene expression profiles



- Limited usage to date one of the biggest findings has been on the ability of certain *E.lenta* strains to metabolize digoxin ([Haiser et al, 2014 Science](#))
- Use has expanded in the last couple of years for viral identifications and novel strain finding



A world of strain biology within individual species

Article

Cell Host & Microbe

The *Prevotella copri* Complex Comprises Four Distinct Clades Underrepresented in Westernized Populations

Graphical Abstract

The figure illustrates the global distribution of *P. copri* prevalence and the presence of four distinct clades (A, B, C, D) across different populations. A world map shows the distribution of *P. copri* prevalence (blue dots for Westernized lifestyle, red triangles for Non-Westernized lifestyle). The Non-Westernized lifestyle is prevalent in 7 countries across 4 continents, while the Westernized lifestyle is prevalent in 22 countries across 3 continents. The study analyzed 6,874 public metagenomes and 272 recently sequenced non-Westernized metagenomes. The *P. copri* complex comprises four distinct clades (A, B, C, D), which are more prevalent in Non-Westernized populations. The Non-Westernized population (red silhouette) contains all four clades (A, B, C, D) in higher proportions compared to the Westernized population (blue silhouette), which contains clade A alone or in combination with clade B.

Authors

Adrian Tett, Kun D. Huang,
Francesco Asnicar,,
Curtis Huttenhower, Frank Maixner,
Nicola Segata

Correspondence

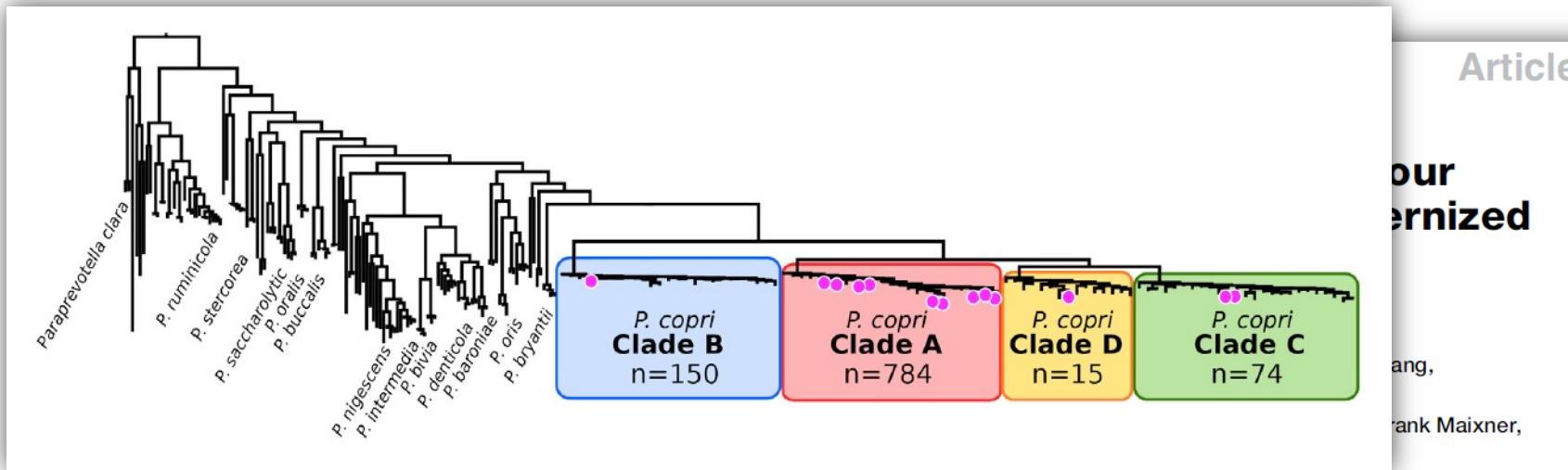
adrianjames.tett@unitn.it (A.T.),
nicola.segata@unitn.it (N.S.)

In Brief

Tett et al. find that the intestinal microbe *Prevotella copri* encompasses four distinct clades constituting the *P. copri* complex. The complex is prevalent in non-Westernized populations where co-presence of all clades is commonly observed within individuals. Analysis of ancient stool samples supports Westernization as contributing to reduced *P. copri* prevalence.



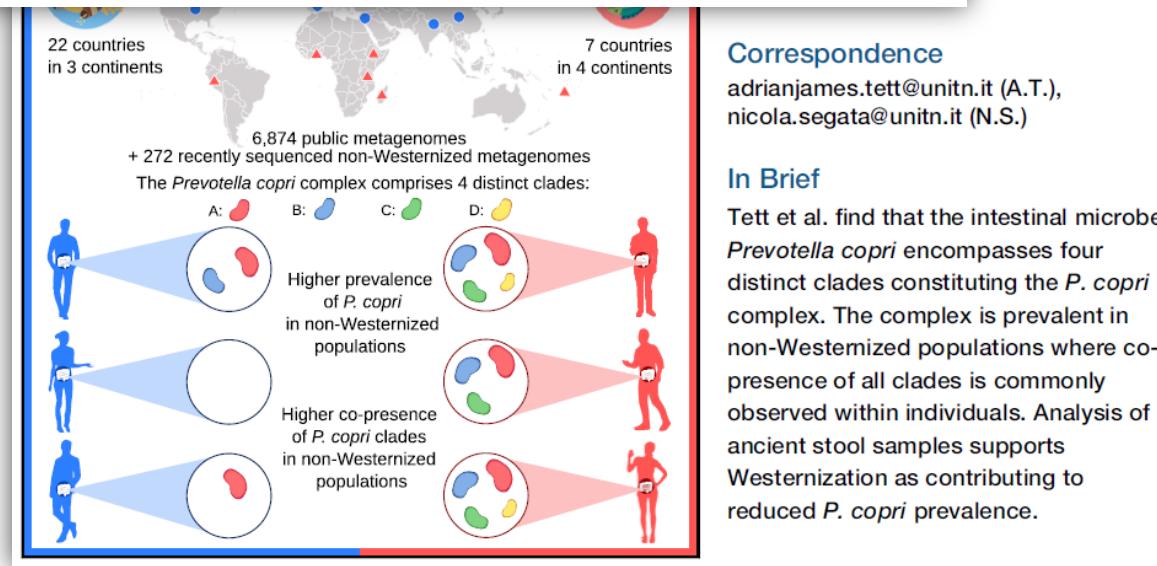
A world of strain biology within individual species



Article

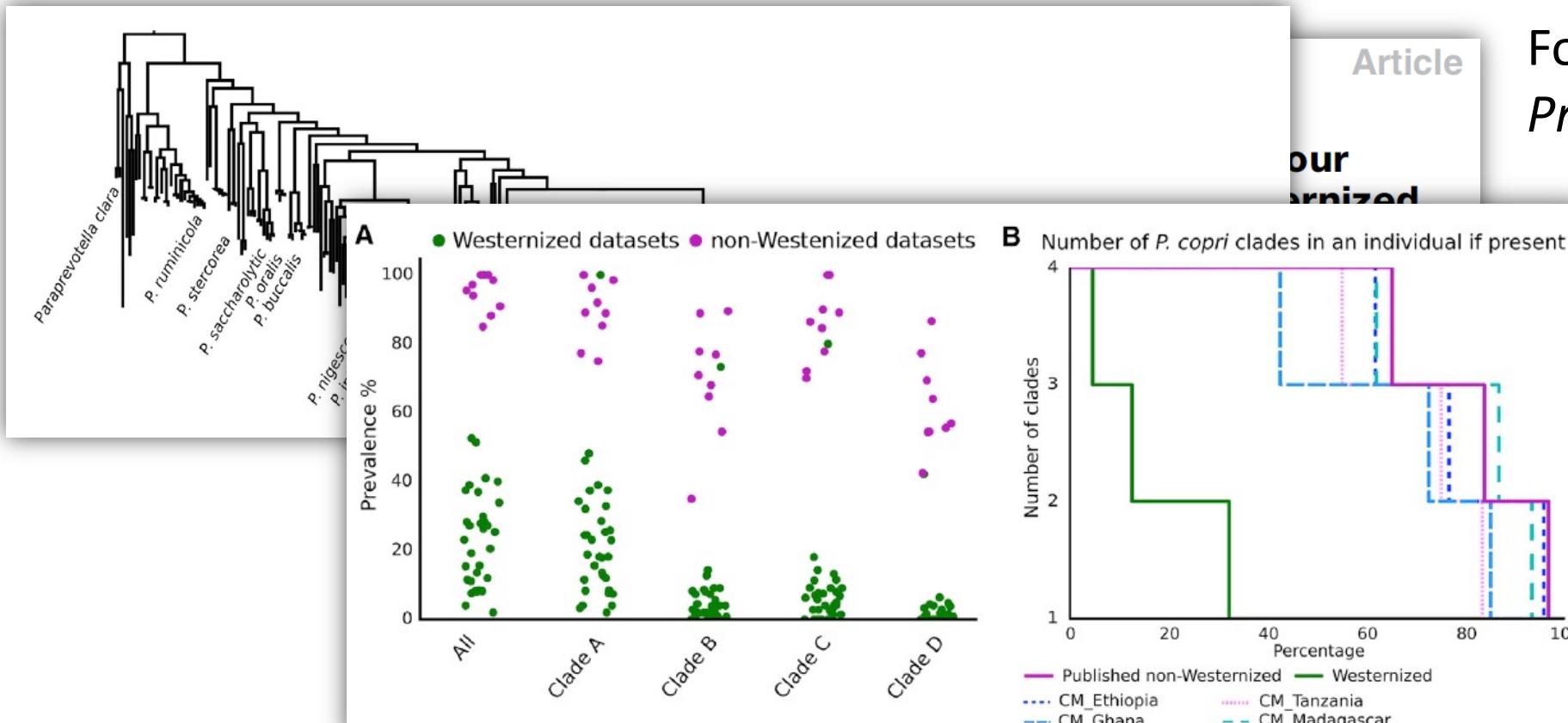
our
ernized

ang,
rank Maixner,

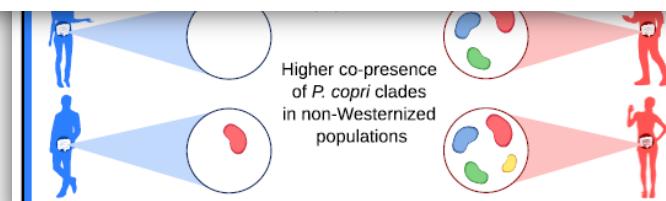


Four clades within the *Prevotella copri* species

A world of strain biology within individual species



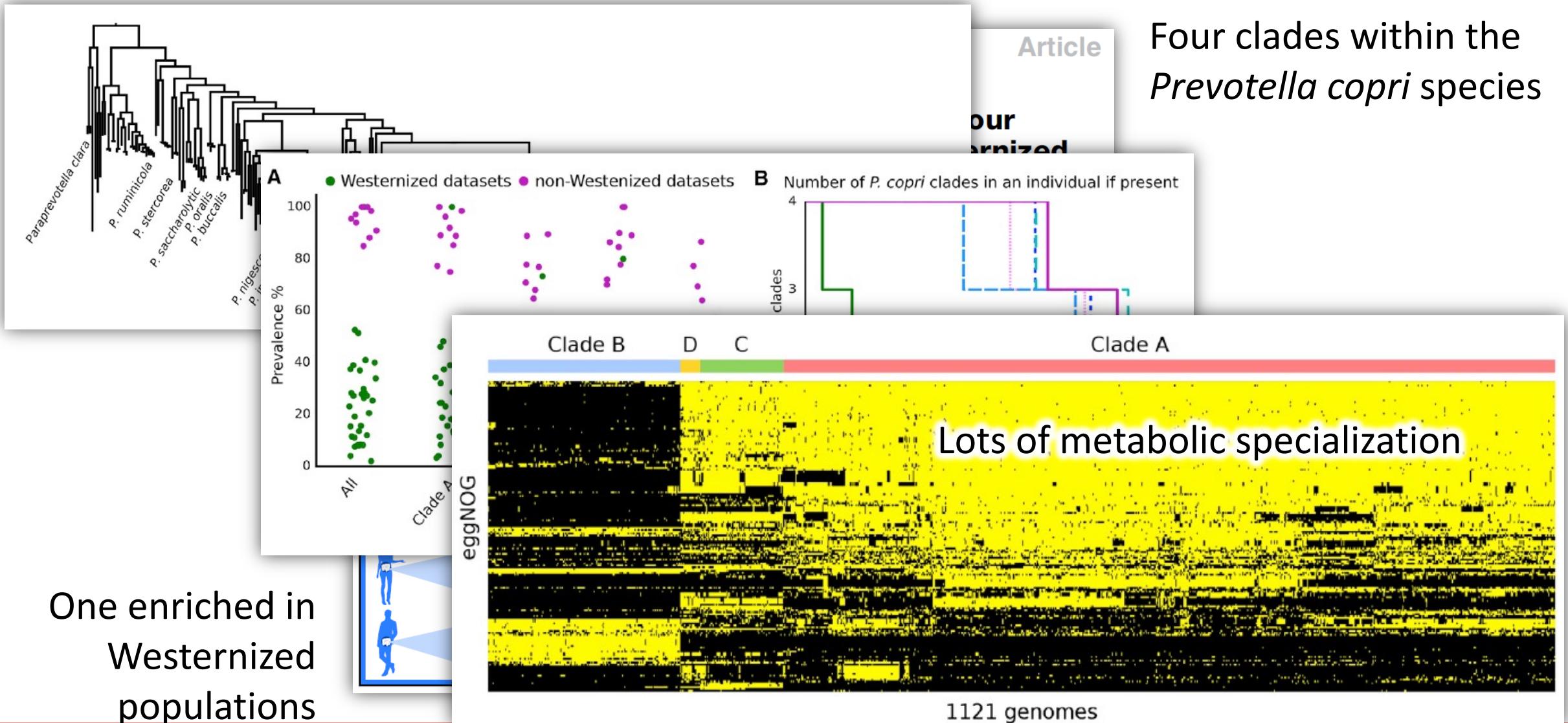
One enriched in
Westernized
populations



presence of all clades is commonly observed within individuals. Analysis of ancient stool samples supports Westernization as contributing to reduced *P. copri* prevalence.



A world of strain biology within individual species

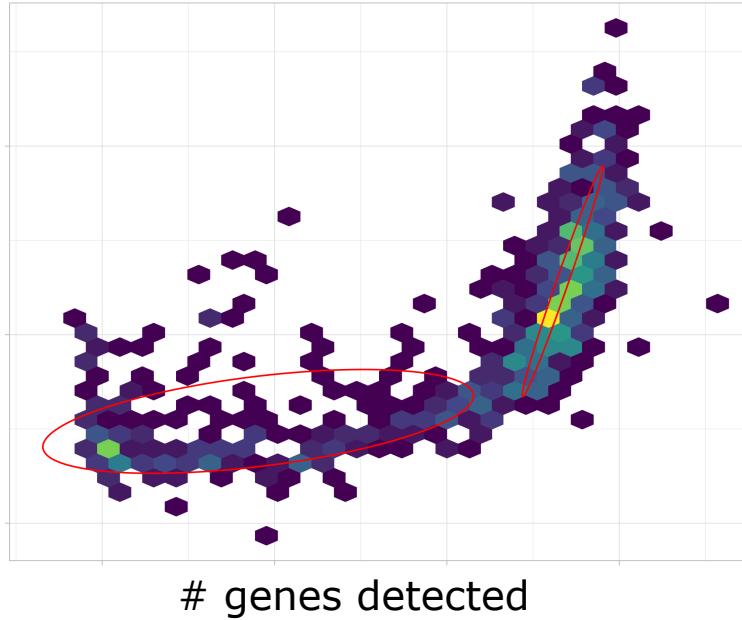


ANPAN: METHODS FOR STRAIN-LEVEL BIOMARKER DISSECTION



Team:
OPTIMISTICC

Evenness of rel. abd.



ANPAN



Andrew
Ghazi

26 July 2022

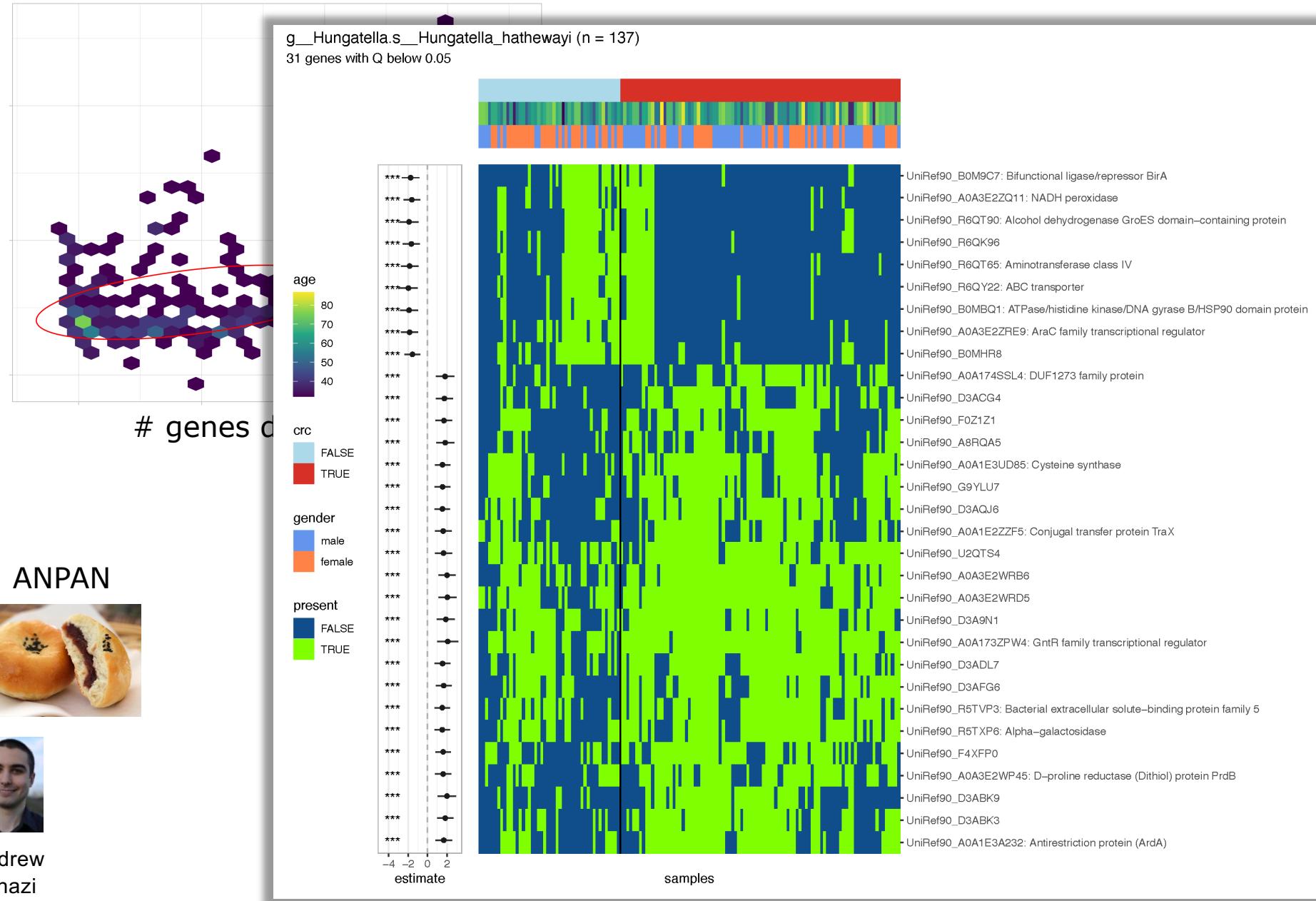
124

ANPAN: METHODS FOR STRAIN-LEVEL BIOMARKER DISSECTION



Team:
OPTIMISTICC

Evenness of rel. abd.



ANPAN



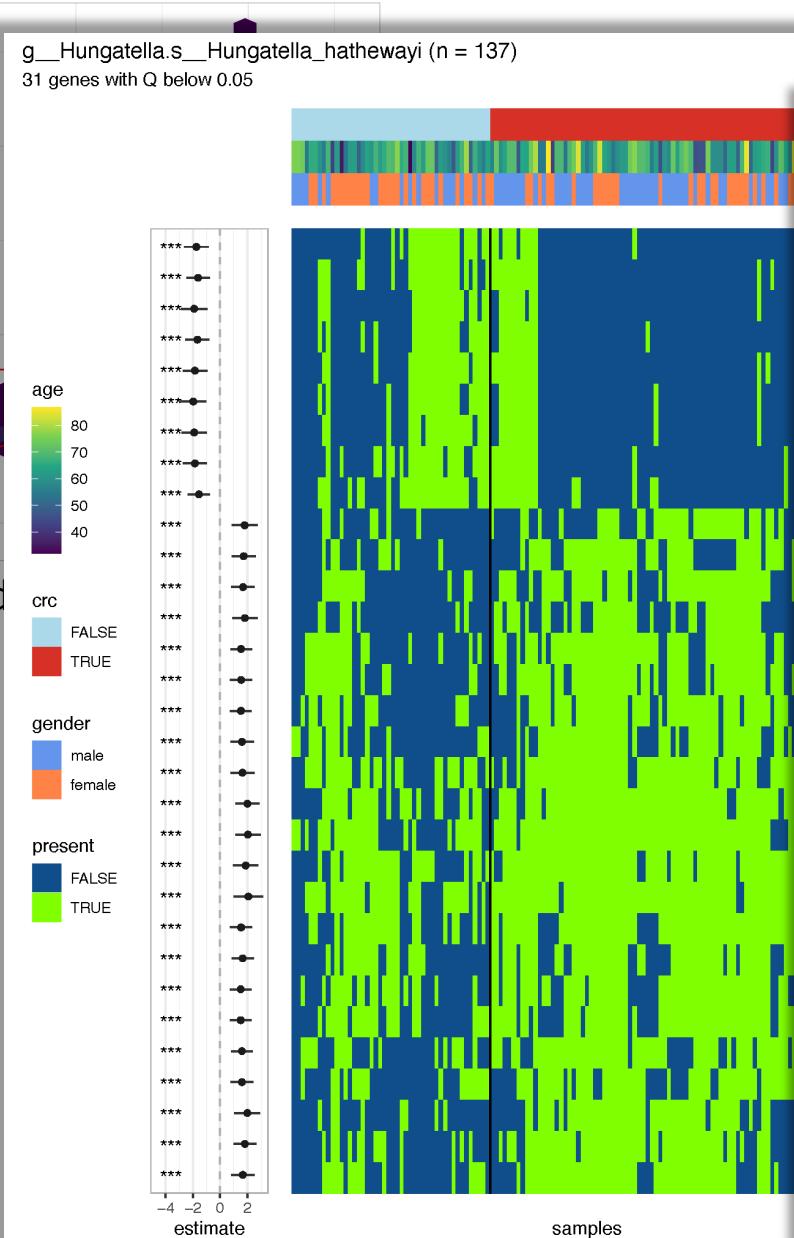
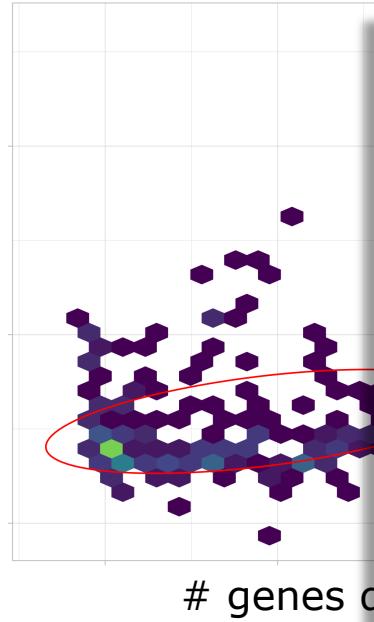
Andrew
Ghazi

ANPAN: METHODS FOR STRAIN-LEVEL BIOMARKER DISSECTION

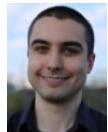


Team:
OPTIMISTICC

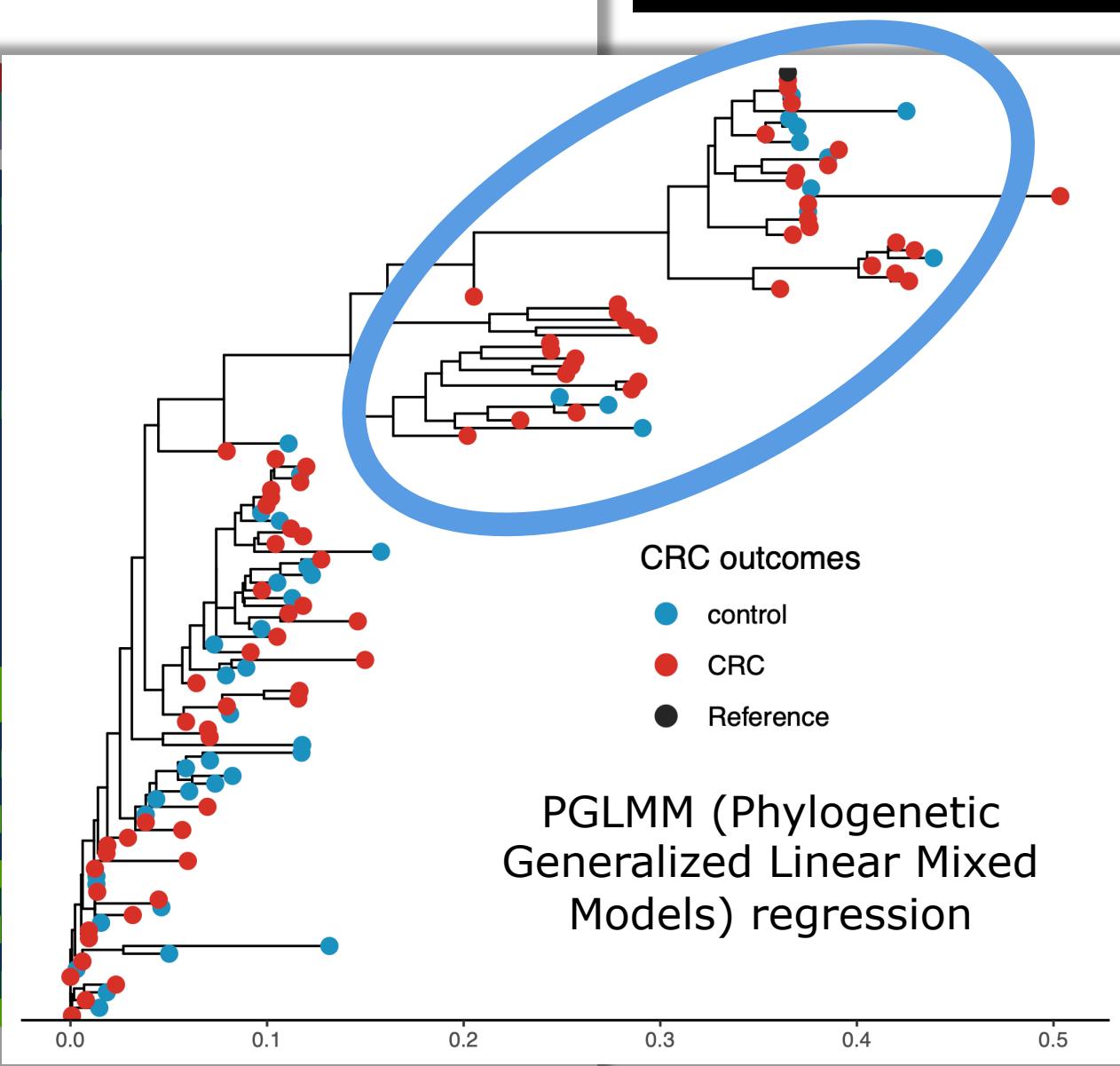
Evenness of rel. abd.



ANPAN



Andrew
Ghazi





Assembly-based approaches



Haha fooled you –
that's tomorrow

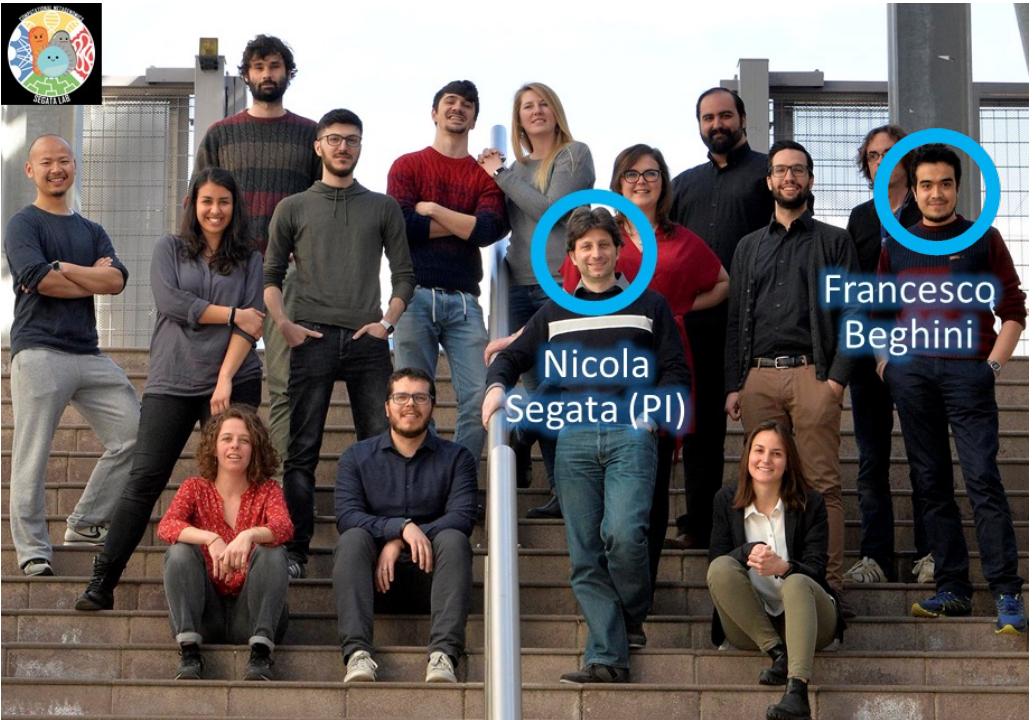
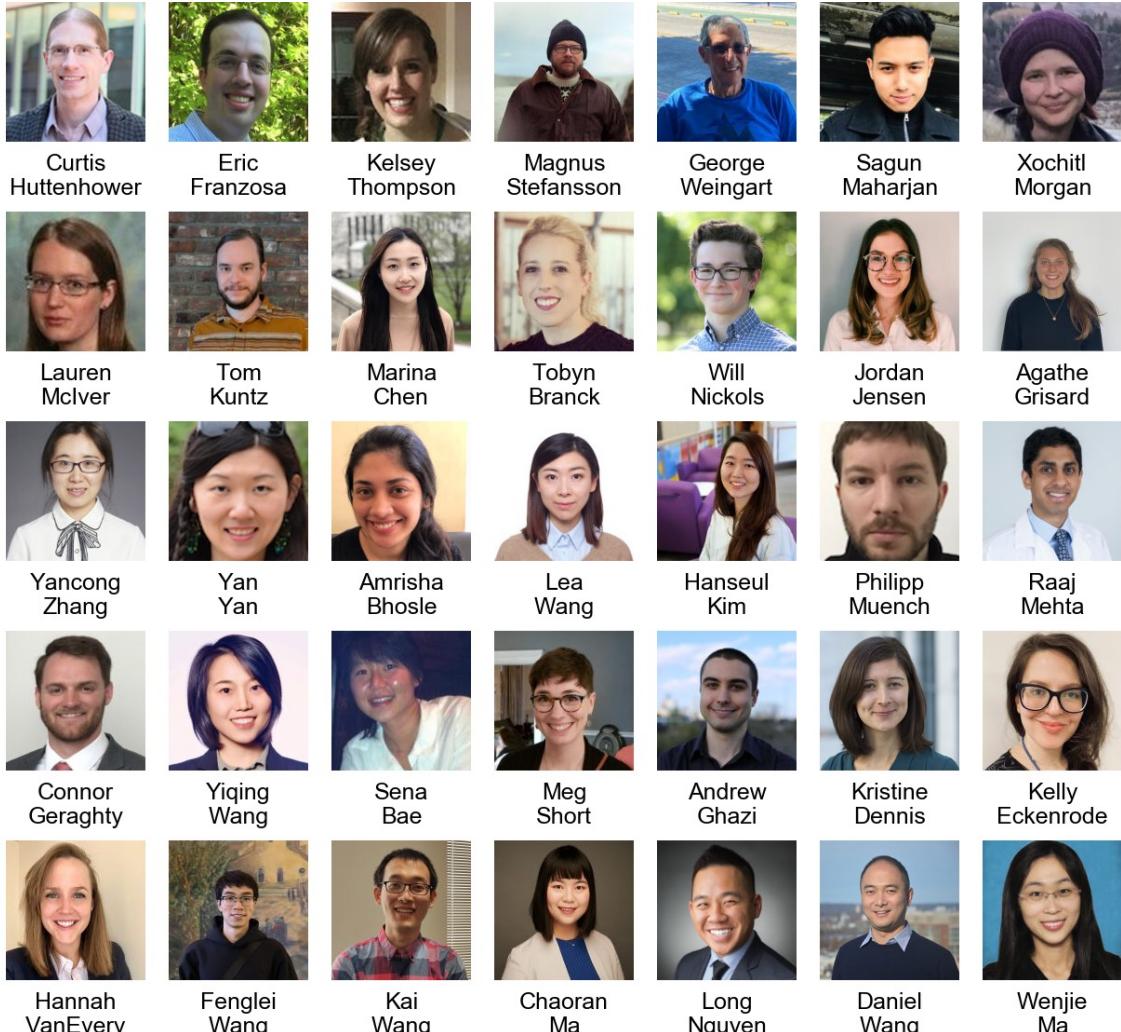


Summary

- What is a microbial strain?
- Why is considering strains important?
- Gene-level strain profiling with PanPhlAn
- SNV-level strain profiling with StrainPhlAn
- Community strain statistics with ANPAN



The bioBakery team(s)



The Laboratory of Computational Metagenomics

University of Trento, Italy

<http://segatalab.cibio.unitn.it>

The Huttenhower Lab

Harvard Chan School of Public Health

<http://huttenhower.sph.harvard.edu>